# ABSTRACT

SILVA, LUCIANO DA COSTA E.  Multiple Trait Multiple Interval Mapping of Quantitative Trait Loci from Inbred Line Crosses. (Under the direction of Zhao-Bang Zeng).

Tremendous progress has been made in recent years on developing statistical methods for mapping quantitative trait loci (QTL) from crosses of inbred lines. Most of the recent research is focused on strategies for mapping multiple QTL and associated model selection procedures and criterion. In Chapter 1, we review the progress of research on QTL mapping on one and multiple trait by maximum likelihood and Bayesian methods.

Although in many instances multiple trait are measured in the same subject, single traits analyses have been the main stream for the purpose of QTL identification. However, single trait analyses do not take advantage of correlation between traits. Multiple trait analysis allows an investigator to assess the pattern of action of QTL on multiple trait, such as, testing the hypothesis of existence of pleiotropic QTL versus the hypothesis of close linked QTL affecting multiple trait, and testing the hypothesis of QTL by environment interaction. In Chapter 2, we proposed a statistical model for mapping multiple QTL affecting multiple trait, the multiple trait multiple interval mapping (MTMIM) model. We also developed a score-based threshold for assessing significance level of QTL effects on multiple trait. Our MTMIM model provides a comprehensive framework for QTL inference in multiple trait, in which the score-based threshold is built in as an essential and elegant tool for computing the significance level of effects of putative QTL in the genome-wide scan, therefore, allowing us to build a set of models containing multiple QTL.

In Chapter 3, we empirically showed that the score-based threshold maintains the false discovery rate within acceptable levels and the multiple trait analysis can bring insights into the analysis of data for the purpose of QTL identification. The analysis of data from an experiment with *Drosophila* showed the potential of our MTMIM model in delivering complementary information regarding the genetic architecture of complex traits, such as, estimating QTL effects on a set of traits simultaneously, testing for the presence of pleiotropic QTL, and estimating the genotypic covariance between traits. A generalized expectation maximization Newton-Raphson (GEM-NR) algorithm for maximizing the likelihood function and estimating parameters in the MTMIM model was compared to the expectation-conditional maximization (ECM) algorithm. Empirical comparison showed that GEM-NR speeded up the convergence of likelihood function considerably when compared to the ECM algorithm, while still delivering stable estimates of parameters.

In Chapter 4, we proposed analytical formulae to predict the length of confidence interval for position of QTL and to predict shape of the LRT around the position of QTL in highly saturate linkage maps and multiple trait analysis using large sample theory. Our results generalize the results of VISSCHER and GODDARD (2004) and they can be used to predict the length of confidence interval for position of QTL with a hypothesized effect on multiple trait, for any given coverage probability. Our analytical formulae can also be used to predict shape of LRT around the position of QTL. Furthermore, we proposed an alternative method for predicting the length of confidence interval for position of QTL, the adjusted method. The adjusted method accounts for the length of the chromosome in which the QTL is located and can deliver more accurately predictions than the method with no adjustments, especially for QTL of low heritability. Our simulation results showed that for sample size of 300 and QTL with heritability levels of 5, 10 and 15%, there are good agreement between lengths of confidence intervals empirically estimated and analytically predicted with the adjusted method.

Multiple Trait Multiple Interval Mapping of
Quantitative Trait Loci from Inbred Line Crosses

by
Luciano Da Costa E Silva

A dissertation submitted to the Graduate Faculty of
North Carolina State University
in partial fullfillment of the
requirements for the Degree of
Doctor of Philosophy

Statistics

Raleigh, North Carolina

2010

APPROVED BY:

_____          _____
Dr. Trudy F.C. Mackay                      Dr. Sujit K. Ghosh


_____          _____
Dr. Zhao-Bang Zeng                         Dr. Howard D. Bondell
Chair of Advisory Committee                Co-Chair of Advisory Committee

# DEDICATION

Aos meus pais Amador e Lurdinha.

To my parents Amador and Lurdinha.

# BIOGRAPHY

Luciano Da Costa E Silva was born in Pitangui, state of Minas Gerais, Brazil. Although he was born in Pitangui, he grew up in the small town of Martinho Campos, also in the state of Minas Gerais. In his childhood he did pretty much all things that a good citizen of Martinho Campos does, playing soccer at the neighbor's lawn, with permission of course, fishing, playing hide and sick, and so on, but he also kept up with his school homework regularly, otherwise he would not be writing this dissertation. As the time went by the small town could not provide more for his parent's dream of his education, then after turning seventeen he was sent to the Central de Ensino e Desenvolvimento Agrário de Florestal, an outstanding technical school located in the city of Florestal, Minas Gerais, where he lived for the next three years until graduation in 1995. During the next one year and seven months he worked as a technician and supervisor for a private company at Brazil's capital, Brasília. In 1998 he began his undergraduate studies towards the bachelor degree in Agronomy in one of the most recognized universities in Brazil, Universidade Federal de Viçosa, located in Viçosa, state of Minas Gerais. In 2003 he finished his undergraduate studies and joined the masters program in Genetics and Plant Breeding in the same university. In 2005 he completed his masters degree and moved to the United States of America to pursue his doctoral degree in Statistics in one of the oldest and most famous Departments of Statistics in the country at North Carolina State University, Raleigh, North Carolina.

# ACKNOWLEDGMENTS

I would like first to express my gratitude to my advisor Dr. Zhao-Bang Zeng, who has always given many bright advices and encouragements not only towards this dissertation but also towards my career. His great generosity and dedication has inspired me.

Drs. Howard D. Bondell, Trudy F.C. Mackay and Sujit K. Ghosh have also provided very helpful feedback when needed. I appreciate your time and dedication.

I would like to thanks my past advisors Drs. Everaldo Gonçalves de Barros, Cosme Damião Cruz, Ney Sussumu Sakiyama and Antônio Alves Pereira, for their valuable advices and help. They inspired me to pursue my doctoral degree and helped me to understand what science is about.

Thanks to Newton Diniz Piovesan, José Manoel Vieira Silva, Daniel Pereira Rocha, Lauro José Moreira Guimarães, and Wagner Luiz Araújo, for their support throughout my life as good friends.

Thanks to Adrian Blue for his hard work at the graduate office in the Department of Statistics. He has always been so generous.

Thanks to the agency Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) and Fulbright, for the financial support, for without this work could not be done so successfully.

My gratitude and deep appreciation to all taxpayers in Brazil and the United States, that financed this work with their hard work.

Thanks to my siblings Admilson, Alisson and Leidiane, for their support throughout my life.

Thanks to my parents Amador and Lurdinha, for their hard work and emotional support given to their children, always encouraging them to pursue higher education.

Finally, my sincerely thanks to my beloved partner Edith Ramos da Conceição Neta. Edith has provided me with the inspiration and emotional support for without this work would have been infinitely harder.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| AIC | Akaike's Information Criterion |
| AICc | Akaike's Information Criterion corrected |
| BC | Backcross |
| BF | Bayes Factor |
| BIC | Bayesian Information Criterion |
| CET | Conditional Empirical Threshold |
| CI | Confidence Interval |
| CIM | Composite Interval Mapping |
| cM | centiMorgans |
| DH | Double Haploid |
| EM | Expectation Maximization |
| ECM | Expectation-Conditional Maximization |
| FDR | False Discovery Rate |
| GEM-NR | Generalized EM-NR |
| GN | Gauss-Newton |

IM              Interval Mapping

LOD             Logarithm of the Odds

LRT             Likelihood Ratio Test

MCMC            Markov Chain Monte Carlo

MIM             Multiple Interval Mapping

MLE             Maximum Likelihood Estimates

MTMIM           Multiple Trait Multiple Interval Mapping

NR              Newton-Raphson

QTL             Quantitative Trait Loci

RET             Residual Empirical Threshold

RIL             Recombinant Inbred Lines

SE              Standard Error

# 1

# Introduction

Many traits that are important to agriculture, human health and evolutionary biology are quantitative in nature, influenced by multiple genes. Efficient and robust identification and mapping onto genomic positions of those genes is a very important goal in quantitative genetics. The availability of genome-wide molecular markers provides the means for us to locate and map those quantitative trait loci (QTL) in a systematic way. Since the publication of LANDER and BOTSTEIN (1989), that first proposed interval mapping method for a genome-wide scan of QTL, many statistical methods have been proposed and developed to map multiple QTL with epistasis in a variety of populations.

In this chapter we review the current status of research on statistical methods for mapping multiple QTL in single and multiple complex traits within the maximum likelihood and Bayesian frameworks.

## 1.1 Inbred populations

The general approach for constructing inbred populations for the purpose of mapping QTL is to choose two homozygous strains (inbred lines) that differ for the traits of interest. Let's denote these parental strains as $P_1$ and $P_2$. Parentals $P_1$ and $P_2$ are

homozygous for all loci. Once the parental strains have been chosen, the next step is to cross them to obtain the filial $F_1$ generation, in which all subjects are genetically identical and heterozygous for all loci at which $P_1$ and $P_2$ differ. Upon having the filial generation and parental strains, and depending on the strategy adopted, many types of experimental inbred populations can be formed, as reviewed in the reminder of this section.

**Intercross $F_2$**

The intercross $F_2$ is obtained by mating subjects of the $F_1$ generation (in some plants, self pollinating). For any locus $Q$ at which the parental strains differ, a subject in the intercross $F_2$ population will have either one of the three genotypes, $QQ$, $Qq$, or $qq$, with expected probability $\frac{1}{4}$, $\frac{1}{2}$, and $\frac{1}{4}$, respectively. For each trait $t$, we parameterize the genotypic values of genotypes $QQ$, $Qq$, $qq$, as $a_t$, $d_t$ and $-a_t$, respectively, where $a_t$ is the additive effect and $d_t$ is the dominance effect of the QTL (FALCONER and MACKAY, 1996). The possibility of sampling all three genotypes in the $F_2$ population allows an investigator to estimate the additive and dominance effects of QTL.

**Recombinant inbred lines**

The first step towards constructing a population of recombinant inbred lines (RIL) is to obtain an $F_2$ generation as describe previously. The second step consists in making repeated mating of relatives (in some plants, self pollinating) for many generations until homozygosis is reached. Plant geneticists usually adopt the single seed decent method starting from the $F_2$ generation to produce the desired RIL. In the single seed decent method, each plant in the $F_2$ population is self pollinated and a single seed from each plant is cultivated to produce a $F_3$ plant. Each $F_3$ plant is self pollinated and a single seed from each $F_3$ plant is cultivated to produce a $F_4$ plant. This process is repeated for many generations, generally until $F_7$, when the RIL population is formed. As the result of mating relatives, each subject in the RIL population is highly homozygous. Therefore, for any locus $Q$ for which the parental lines differ each

subject has almost certainly either genotype $QQ$ or $qq$. Therefore, an investigator can obtain information only about the additive effect of QTL in RIL populations.

### Double haploid

The subjects of a double haploid (DH) population are obtained from regeneration *in vitro* of immature pollen grains with their genetic material duplicated. In the technique, pollen grains are harvested from adult $F_1$ plants and chemically induced to duplicate their chromosomes without further cell division. These pollen grains with duplicated chromosomes are cultivated *in vitro* to further develop into plant siblings. Similar to RIL populations, subjects in DH populations have only homozygous genotypes. Therefore, only the additive effect of QTL can be studied. The use of DH population is very common among rice breeders.

### Backcross

The backcross (BC) population is obtained with the crossing of $F_1$ subjects to one parental strains. For any locus $Q$ in which the parental strains differ, any subject in the BC population has either genotype $QQ$ or $Qq$, if the genotype of parental strain crossed to $F_1$ is $QQ$. Because only the homozygous genotype (QQ) of a segregating locus is sampled in the BC population, estimation of additive and dominance effects of a QTL is impossible.

Although many population structures can be used for studying the effects of QTL, we have chosen to discuss the statistical methods throughout this review with the use of a BC population, in which the filial generation $F_1$ of two inbred lines have been backcrossed to the parental line that has genotype $QQ$ for a QTL.

Suppose that parental strains $P_1$ and $P_2$ have genotypes $m_L m_L q q m_R m_R$ and $M_L M_L Q Q M_R M_R$, respectively, where, $M$ and $m$ are genetic marker alleles of marker loci on the left (L) and right (R) of QTL, and $Q$ and $q$ are alleles of QTL. Any offspring $F_1$ from the cross between $P_1$ and $P_2$ would have genotype $M_L m_L Q q M_R m_R$, and any subject of a population obtained from the cross between $F_1$ and $P_2$ would have one of the eight genotypes shown in Table 1.1. The expected frequency of each

genotype is also shown in Table 1.1. We assume: (1) the three loci are linked (i.e., the loci are relatively close to each other on a chromosome); (2) the order of loci in the chromosome is $M_L Q M_R$; (3) the recombination frequencies between $M_L$ and $Q$, $Q$ and $M_R$, and $M_L$ and $M_R$ are, $r_{LQ}$, $r_{QR}$, and $r_{LR}$, respectively.

Table 1.1: Joint probabilities of genotypes for any subject originated from the backcross of $F_1$ ($M_L m_L Q q M_R m_R$) and $P_2$ ($M_L M_L Q Q M_R M_R$), assuming the loci are linked.

| Genotypes[1] | Probability[2] |
|---|---|
| $M_L M_L Q Q M_R M_R$ | $\frac{1}{2}[(1 - r_{LQ})(1 - r_{QR}) - r_{LQ} r_{QR}(1 - c)]$ |
| $M_L M_L Q Q M_R m_R$ | $\frac{1}{2}[(1 - c r_{LQ}) r_{QR}]$ |
| $M_L M_L Q q M_R M_R$ | $\frac{1}{2}[c r_{LQ} r_{QR}]$ |
| $M_L M_L Q q M_R m_R$ | $\frac{1}{2}[r_{LQ}(1 - c r_{QR})]$ |
| $M_L m_L Q Q M_R M_R$ | $\frac{1}{2}[r_{LQ}(1 - c r_{QR})]$ |
| $M_L m_L Q Q M_R m_R$ | $\frac{1}{2}[c r_{LQ} r_{QR}]$ |
| $M_L m_L Q q M_R M_R$ | $\frac{1}{2}[(1 - c r_{LQ}) r_{QR}]$ |
| $M_L m_L Q q M_R m_R$ | $\frac{1}{2}[(1 - r_{LQ})(1 - r_{QR}) - r_{LQ} r_{QR}(1 - c)]$ |

[1] The order of the markers and QTL is assumed to be $M_L Q M_R$, and the recombination frequencies between $M_L Q$, $Q M_R$, and $M_L M_R$ are, respectively, $r_{LQ}$, $r_{QR}$, and $r_{LR}$.

[2] $c = \frac{OFDR}{EFDR} = \frac{OFDR}{r_{LQ} r_{QR}}$, where OFDR and EFDR stand for observed and expected frequency of double recombinants, respectively.

Unfortunately, the assessment of probabilities in Table 1.1 is not possible because genotypes of QTL are unobserved. Because the only information available to us is the genotypes of genetic markers flanking QTL, the general strategy is to compute, for each subject, the conditional probabilities of QTL genotypes $QQ$ and $Qq$, given the two markers $M_L$ (on the left) and $M_R$ (on the right) flanking the QTL (Table 1.2). In cases where missing information on the markers flanking the QTL is present,

one may use the information on the closest neighboring markers through a Markov Chain (JIANG and ZENG, 1997).

Table 1.2: Conditional probabilities of QTL genotypes $QQ$ and $Qq$, for any subject of a BC population, given the genotypes of markers $M_L$ (on the left) and $M_R$ (on the right) flanking the QTL. The assumption of complete cross-over interference ($c = 0$) within a marker interval was made in this table.

| Marker genotypes | $P(QQ|M_L, M_R, r_{LR}, r_{LQ})$ | $P(Qq|M_L, M_R, r_{LR}, r_{LQ})$ |
|---|---|---|
| $M_L M_L M_R M_R$ | $1$ | $0$ |
| $M_L M_L M_R m_R$ | $\frac{r_{LQ}}{r_{LR}}$ | $1 - \frac{r_{LQ}}{r_{LR}}$ |
| $M_L m_L M_R M_R$ | $1 - \frac{r_{LQ}}{r_{LR}}$ | $\frac{r_{LQ}}{r_{LR}}$ |
| $M_L m_L M_R m_R$ | $0$ | $1$ |

$r_{LR}$ is the recombination frequency between markers $M_L$ and $M_R$, and $r_{LQ}$ is the recombination frequency between locus $M_L$ and QTL ($0 \leq r_{LQ} \leq r_{LR}$).

## 1.2  Outbred populations

In the previous section we reviewed some commonly inbred populations used for mapping QTL. For some populations, for instance, livestock and humans, developing and crossing inbred lines are not feasible, however. The information available for such populations is collected in terms of pedigrees over multiple generations.

In livestock, many family structures have been described in the literature (VAN DER BEEK *et al.*, 1995), such as, half-sibs and full-sibs. In half-sibs, one of the parents, generally the male, has many unrelated mates and each mate has one offspring, in which the traits are measured. Genotypes are obtained for the common parent and offsprings. In full-sib, a pair of parents has many offsprings, in which the trait values are obtained. Marker genotypes are obtained for each offspring and its pair of parents.

Since the distances between QTL and markers in a moderate resolution linkage map (markers ten centiMorgans (cM) apart from each other) often exceeds 1 cM, recombinations between QTL and markers over generations prevent the detection of any linkage disequilibrium resulting from a not recent event in the history of the population. Therefore, linkage disequilibrium information is only available within parents, and hence QTL effects must be estimated within parents (HOSCHELE, 2007).

Complexity of structure of populations dictates the complexity of the statistical model for mapping QTL. When large individual families are available, for instance, in cattle a single sire may have hundreds or even thousands of offsprings, simple adaptations of methods for inbred populations are feasible (KNOTT *et al.*, 1996). In contrast, variance-component methods is advocated (HOESCHELE *et al.*, 1997) to analyze multiple families of small sizes and possibly with genetic ties between families.

Statistical methods for mapping QTL in outbred populations would deserver, on their own, a rigorous review chapter. However, we omit further details in this dissertation to dedicate more efforts towards inbred populations, which are used for modeling purpose throughout this dissertation.

## 1.3   Mapping multiple QTL on single trait

The identification of QTL effects associated with genetic markers dates back to SAX (1923). While the first method simply tested for differences between means of a phenotype associate with marker genotypes by $t$-statistic (SOLLER and BRODY, 1976; REBAI *et al.*, 1995; LIU, 1998; WU *et al.*, 2007), more elaborated methods such as interval mapping (IM) (LANDER and BOTSTEIN, 1989) use information on more than one genetic marker at a time, therefore, delivering more power in identifying QTL and in estimating their effects. Because the IM assumes a single QTL in the model, the effects of unaccounted QTL remains in the model residual sum of squares. Therefore, IM does not take advantage of information on multiple markers to increase the power of the test statistic and consequently its ability in identifying a QTL in the presence

of other QTL. JANSEN (1993) and ZENG (1993, 1994), independently, proposed a method based on multiple regression that fits both the effect of a QTL as well as the effects of covariates, which are a subset of selected genetic markers. The method is a combination of interval mapping and multiple regression. The interval mapping is used to fit a linear model at every position for the genome scan, and the multiple regression is used for fitting covariates to control linked and unlinked QTL effects and reduce the model residual. When well-chosen, the use of covariates reduces the model residual sum of squares, enhancing the power of identifying QTL. The multiple regression method proposed by Zeng is known in the literature as composite interval mapping (CIM).

In this section, we review statistical methods for identifying QTL with significant effects on the distribution of a single quantitative trait. There is a vast literature on both maximum likelihood and Bayesian methods for mapping QTL and we will comment them along the way. In what follows, for any matrix $\boldsymbol{A}$, its transpose is denoted by $\boldsymbol{A}'$, its inverse by $\boldsymbol{A}^{-1}$, its $u^{th}$ row by $\boldsymbol{A}_{[u,\cdot]}$, its $v^{th}$ column by $\boldsymbol{A}_{[\cdot,v]}$, and its element in row $u$ and column $v$ by $\boldsymbol{A}_{[u,v]}$.

## 1.3.1 Multiple interval mapping: maximum likelihood method

In this section, we describe the multiple interval mapping (MIM) (KAO and ZENG, 1997; KAO *et al.*, 1999; ZENG *et al.*, 1999) with details. MIM is similar to CIM in concept. The rationale of CIM is to fit in the model genetic markers closely linked to other QTL across the genome as covariates when searching for a specific QTL. The rationale of MIM is to fit in the model estimated positions of other QTL rather than their closely linked genetic markers, therefore, delivering more power as well as better parameter estimates.

**Model**: The statistical model for single trait multiple QTL inference on BC population is a linear model in which, the trait value $y_i$ of the $i^{th}$ subject ($i = 1, 2, \cdots, n$) is regressed on the explanatory variables $x_{ir}$ ($r = 1, 2, \cdots, m$), which are defined according to the Cockerham genetic model (KAO *et al.*, 1999; KAO and ZENG, 2002;

ZENG *et al.*, 2005). For each subject $i$ in our BC population defined previously, $x_{ir}$ takes value $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether the QTL $r$ has genotype $QQ$ or $Qq$, respectively. The coefficients of $x_{ir}$, $\beta_r$, are called the effect of the $r^{th}$ QTL. The linear model also includes an intercept $\mu$, the epistatic effects $w_{rl}$ between QTL $r$ and $l$ for a subset $p$ of all pairwise interactions, and the residue $e_i$, which is assumed to be independent and identically distributed according to a normal distribution with mean zero and variance $\sigma_e^2$. The linear model is then:

$$y_i = \mu + \sum_{r=1}^{m} \beta_r x_{ir} + \sum_{r<l}^{p} w_{rl} x_{ir} x_{il} + e_i \tag{1.1}$$

Let $\boldsymbol{y}$ be the $n$ by 1 vector of all observations, $\boldsymbol{X}$ be the $n$ by $m+p$ incidence matrix, $\boldsymbol{e}$ be the $n$ by 1 vector of residuals, $\boldsymbol{1}$ be an $n$ by 1 vector of ones, and let $\boldsymbol{\beta} = (\beta_1, \beta_2, \cdots, \beta_m, w_1, w_2, \cdots, w_p)'$ be the vector of all main and epistatic effects and $\boldsymbol{\theta} = (\boldsymbol{\beta}', \mu, \sigma_e^2)$. Then the statistical model in matrix form would look like $\boldsymbol{y} = \boldsymbol{1}\mu + \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{e}$.

**Likelihood**: Let $\boldsymbol{M}$ be an $n$ by $K$ matrix of known marker genotypes and let $\boldsymbol{\mathcal{R}} = (r_{12}, r_{23}, \cdots, r_{K-1,K})$ be a vector of known recombination frequencies (or distances) between markers, where $r_{k,k+1}$ is the recombination frequency (or distance) between markers $M_k$ and $M_{k+1}$. In the BC population defined previously there are two possible genotypes for each QTL, $QQ$ and $Qq$. Therefore, if there are $m$ loci affecting a trait, there are $2^m$ possible genotypes for any subject $i$, genotypes of the form $G_j = Q_1 Q_2 \cdots Q_m$, where $Q_r \in \{QQ, Qq\}$, $r = 1, 2, \cdots, m$ and $j = 1, 2, \cdots, 2^m$. Assuming $m$ main and $p$ epistatic effects in model (1.1), we define an $2^m$ by $s = m + p$ matrix $\boldsymbol{D}$ of coded genotypes according to the Cockerham genetic model (KAO and ZENG, 1997, 2002; ZENG *et al.*, 2005). In the matrix $\boldsymbol{D}$ each column $b$ corresponds to the $b^{th}$ effect parameter ($b = 1, 2, \cdots, s$) and each row $j$ of $\boldsymbol{D}$, $\boldsymbol{D}_{[j,\cdot]}$, represents a coded genotype $G_j$. If $b \leq m$, $\boldsymbol{D}_{[j,b]} = x_r$, otherwise $\boldsymbol{D}_{[j,b]} = x_r * x_l$, where $x_u$ ($u = r$ and $u = l$) is either $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether the genotype of QTL $Q_u$ in $G_j$ is $QQ$ or $Qq$, respectively.

In order to search the entire genome for significant effects of QTL, the genome is partitioned into $H$ loci, usually at 1-cM grid. This partition is denoted by $\boldsymbol{\zeta}$. The set of positions of $m$ putative QTL in model (1.1), $\boldsymbol{\lambda} = \{\lambda_1, \lambda_2, \cdots, \lambda_m\}$, is assumed to be a fixed and known subset of $\boldsymbol{\zeta}$. For any subject $i$, let $\boldsymbol{M}_{[i,\cdot]}$ be the genotypic information of the markers flanking the $m$ QTL, and $M_{i,L}^r$ and $M_{i,R}^r$ be the flanking markers on left and right of QTL $Q_r$, respectively. Then, assuming no cross-over interference between marker intervals and no more than one QTL existing within a marker interval, the probability of any genotype $G_j$ conditional on the genotypes of the markers flanking the $m$ QTL is $P(G_j|\boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\mathcal{R}}, \boldsymbol{\lambda}) = \prod_{r=1}^{m} P(Q_r|M_{i,L}^r, M_{i,R}^r, \lambda_r)$, where the conditional probabilities on the right hand side of the equation can be found in JIANG and ZENG (1997) and KAO and ZENG (1997). In Table 1.2, we show how to estimate these probabilities for a BC population. The recombination frequency between QTL $Q_r$ and its left marker ($r_{LQ_r}$) has a one-to-one correspondence with the position $\lambda_r$. Conditional probabilities of two QTL lying within a single interval is shown in Table A.2 of Appendix A.

If $G_j$ genotype of the $i^{th}$ subject were known, then individual likelihood $L_i$ would be $L_i\left(\boldsymbol{\theta} \mid y_i\right) = \phi\left(y_i|\mu + \boldsymbol{D}_{[j,\cdot]}\boldsymbol{\beta}, \sigma_e^2\right)$, where $\phi(z|\mu_0, \sigma_0^2)$ is the probability density function of a normal random variable $z$ with mean $\mu_0$ and variance $\sigma_0^2$. However, since genotype $G_j$ is unknown, then the individual likelihood ($L_i$), assuming $m$ QTL with positions defined in $\boldsymbol{\lambda}$, is a mixture of $2^m$ normal distributions with homogeneous variance, different means, and mixture probabilities $p_{ij} = P(G_j|\boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\mathcal{R}}, \boldsymbol{\lambda})$:

$$L_i\left(\boldsymbol{\theta} \mid y_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right) = \sum_{j=1}^{2^m} p_{ij}\phi\left(y_i|\mu + \boldsymbol{D}_{[j,\cdot]}\boldsymbol{\beta}, \sigma_e^2\right) \tag{1.2}$$

and the overall likelihood ($L$) is $L\left(\boldsymbol{\theta} \mid \boldsymbol{y}, \boldsymbol{M}, \boldsymbol{\lambda}\right) = \prod_{i=1}^{n} L_i\left(\boldsymbol{\theta} \mid y_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right)$. The maximum likelihood estimates of all parameters in model (1.1) are obtained after maximizing $L$ with some technique, such as, Newton-Rapson, quasi-likelihood, expectation-maximization (EM) algorithm or some of its variants (MCLACHLAN and KRISHNAN,

1996; KAO *et al.*, 1999).

A drawback of model (1.1) when applied to the analysis of data from a BC population is the lack of separate estimates of the additive and the dominance effects of a QTL. An alternative that enables estimates of both additive and dominance effects is to use two backcrosses (a backcross of filial generation $F_1$ to the parental line with QTL genotype $QQ$ and other backcross of $F_1$ to the parental line with genotype $qq$) or NC design III (GARCIA *et al.*, 2008). The method of Garcia also has an additional advantage of allowing for inference on the genetic bases of heterosis.

In what follows we define the logarithm likelihood ratio test (LRT), the logarithm of the odds (LOD) and the score statistic. Then we formalize the concept of genome-wide scan, review methods of threshold computation for a reliable detection of QTL, and end with a review of methods for model selection.

**Hypothesis testing**: In general, let $\boldsymbol{\theta}$ be the vector of all parameters in model (1.1), let $\boldsymbol{\theta}_1 \in \boldsymbol{\theta}$ be a vector of length $c$ for which one wants to test the hypotheses $H_0 : \boldsymbol{\theta}_1 = \boldsymbol{\theta}_{10}$ versus $H_1 : \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_{10}$, where $\boldsymbol{\theta}_{10}$ is a vector of real numbers, let $\boldsymbol{\eta}$ be the vector of nuisance parameters, and let $L(\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\eta})|\boldsymbol{y})$ be likelihood function of the data $\boldsymbol{y}$. In the remainder of this section we describe LOD, LRT and score statistic for testing $H_0$ vs $H_1$ in a general framework, and then for testing QTL effects in genome-wide scan. The *logarithm of the odds* (LOD) is:

$$LOD = \log_{10}\left(\frac{sup_{\boldsymbol{\theta} \in H_0 \cup H_1} L(\boldsymbol{\theta}|\boldsymbol{y})}{sup_{\boldsymbol{\theta} \in H_0} L(\boldsymbol{\theta}|\boldsymbol{y})}\right)$$
$$= \log_{10}\left(\frac{L(\hat{\boldsymbol{\theta}}|\boldsymbol{y})}{L(\tilde{\boldsymbol{\theta}}|\boldsymbol{y})}\right)$$

where, $sup$ stands for sumpremum, and $\hat{\boldsymbol{\theta}}$ and $\tilde{\boldsymbol{\theta}}$ are the maximum likelihood estimates (MLE) under $H_1$ (full model) and $H_0$ (reduced model), respectively.

The *logarithm likelihood ratio test* (LRT) is:

$$LRT = -2\log_e \left( \frac{sup_{\boldsymbol{\theta} \in H_0} L(\boldsymbol{\theta}|\boldsymbol{y})}{sup_{\boldsymbol{\theta} \in H_0 \cup H_1} L(\boldsymbol{\theta}|\boldsymbol{y})} \right)$$

$$= -2\log_e \left( \frac{L(\tilde{\boldsymbol{\theta}}|\boldsymbol{y})}{L(\hat{\boldsymbol{\theta}}|\boldsymbol{y})} \right)$$

The *score statistic* has been in the literature for quite a while since it was proposed by RAO (1948), but it had not been used very often in the QTL mapping literature, until ZOU *et al.* (2004) proposed a score-based threshold for assessing genome-wide significance of QTL effects. In general, let $\ell_i = \log_e (L_i(\boldsymbol{\theta}|y_i))$ and $\ell = \log_e (L(\boldsymbol{\theta}|\boldsymbol{y}))$ be the natural logarithm of the individual and overall likelihoods, respectively.

The score statistic to test $H_0$ vs $H_1$ can be written as $S = \hat{U}'\hat{V}^{-1}\hat{U}$ (ZOU *et al.*, 2004; COX and HINKLEY, 1974), where $\hat{U} = \sum_{i=1}^{n} \hat{U}_i$, $\hat{V} = \sum_{i=1}^{n} \hat{U}_i\hat{U}_i'$, and $\hat{U}_i$ is:

$$\hat{U}_i = \frac{\partial \ell_i(\boldsymbol{\theta}_1, \boldsymbol{\eta})}{\partial \boldsymbol{\theta}_1'} \bigg|_{(\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{10}, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})}$$
$$- \frac{\partial \ell(\boldsymbol{\theta}_1, \boldsymbol{\eta})}{\partial \boldsymbol{\theta}_1 \partial \boldsymbol{\eta}'} \bigg|_{(\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{10}, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})} \left( \frac{\partial \ell(\boldsymbol{\theta}_1, \boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}'} \bigg|_{(\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{10}, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})} \right)^{-1} \frac{\partial \ell_i(\boldsymbol{\theta}_1, \boldsymbol{\eta})}{\partial \boldsymbol{\eta}'} \bigg|_{(\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{10}, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})}$$

$$(1.3)$$

where, $\tilde{\boldsymbol{\eta}}$ is the MLE of $\boldsymbol{\eta}$ under $H_0$.

Under the large sample assumption the LRT and score statistics follows a $\chi^2$ distribution with $c$ degrees of freedom. There is an one-to-one mathematical relationship between LRT and LOD statistics, which is $LRT = 2log_e(10)LOD \approx 4.6LOD$.

**Genome-wide scan**: In genome-wide scan a putative QTL is assumed at every position $\lambda \in \boldsymbol{\zeta}$ and its effect significance (main or epistatic effect) is tested against the null of no effect. For instance, assume a model with $m-1$ main effects and $p$ epistatic effects and we are scanning for a putative $m^{th}$ QTL. Let $l = \lambda$ denotes the testing position of the putative QTL coming into the model and $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \cdots, \lambda_{(m-1)}, l)$ be the current positions of all QTL. Let $\theta_m = \beta_m$ represents the main effects for the new

QTL coming into the model, and let $\boldsymbol{\theta} = (\theta_1, \theta_2, \cdots, \theta_{m-1}, \theta_m, \theta_{m+1}, \cdots, \theta_s, \mu, \sigma_e^2)'$ be the column vector of all parameters, where $\theta_b = \beta_b$ for $1 \leq b \leq m$ and $\theta_b = w_b$ for $m < b \leq s = m + p$. Let $\boldsymbol{\eta} = (\theta_1, \theta_2, \cdots, \theta_{m-1}, \theta_{m+1}, \cdots, \theta_s, \mu, \sigma_e^2)'$ be the column vector of nuisance parameters. Then for every position $l$ either LOD, LRT or score statistics may be used to assess the strength of the hypotheses $H_0 : \theta_m = 0$ versus $H_1 : \theta_m \neq 0$:

$$LOD(l) = \log_{10}\left(\frac{L(\hat{\boldsymbol{\theta}}|\boldsymbol{y},\boldsymbol{M},\boldsymbol{\lambda})}{L(\tilde{\boldsymbol{\theta}}|\boldsymbol{y},\boldsymbol{M},\boldsymbol{\lambda})}\right), \quad LRT(l) = -2\log_e\left(\frac{L(\tilde{\boldsymbol{\theta}}|\boldsymbol{y},\boldsymbol{M},\boldsymbol{\lambda})}{L(\hat{\boldsymbol{\theta}}|\boldsymbol{y},\boldsymbol{M},\boldsymbol{\lambda})}\right), \quad S(l) = \hat{U}'\hat{V}^{-1}\hat{U}$$

where, $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \hat{\theta}_2, \cdots, \hat{\theta}_s, \hat{\mu}, \hat{\sigma}_e^2)$ and $\tilde{\boldsymbol{\theta}} = (\tilde{\theta}_1, \tilde{\theta}_2, \cdots, \tilde{\theta}_{m-1}, 0, \tilde{\theta}_{m+1}, \tilde{\theta}_s, \tilde{\mu}, \tilde{\sigma}_e^2)$ are the MLE under $H_1$ and $H_0$, respectively.

For every position $l$, LRT and score statistic, under some regular conditions, are asymptotically equivalent. Therefore, under a reasonable sample size, the chi-squared distribution $(\chi^2)$ can be used to obtain the appropriate threshold for a given level of significance for a fixed testing position $l$. However, in the MIM genome-wide scan the hypothesis is tested many times, therefore, some correction for the test multiplicity is required. Next, we describe two approaches to tackle this issue, the permutation test and the score-based threshold.

**Genome-wide threshold**: In the genome-wide search for identifying QTL effects associated with certain markers, an investigator wants to increase the chances of finding significant QTL throughout the genome of an organism. An obvious strategy is to maximize the number of screened markers up to the budget and labor work available. If all markers in a large set of markers are tested for the hypothesis of presence of a QTL effect against the null of no effect, then multiple testing is an issue that must be dealt with. The consequence of multiplicity of tests could be even worse in the IM and MIM methods because not only all markers but also many other positions between markers are tested for the presence of a putative QTL. Therefore, some correction is required for controlling the rate of rejecting the null of no effect when it is true (false positive discovery). Theoretical threshold values were proposed

by LANDER and BOTSTEIN (1989) and REBAI *et al.* (1994). The threshold of LANDER and BOTSTEIN (1989) is derived with assumptions of large sample size, infinitely many genetic markers and normal distribution of quantitative trait. REBAI *et al.* (1994)'s threshold is an approximation of DAVIES (1977) and DAVIES (1987)'s bound and requires integral evaluation for each marker interval, which in more complex models must be substituted with numerical integration. Neither one of the two methods of threshold computation is popular, either because of unrealistic assumptions or practical difficulties regarding their implementations.

Recently, ZOU and ZENG (2008) have provided a good review of many strategies for controlling for false positive discovery in QTL analyses, such as, false discovery rate (FDR) (BENJAMINI and HOCHBERG, 1995), Bonferroni correction, permutation (CHURCHILL and DOERGE, 1994), and score-based threshold (ZOU *et al.*, 2004). The low number of discoveries in the genome-wide search limits the practical implementation of the FDR in traditional QTL experiments. As for the Bonferroni correction, it is known that it can be very conservative when there is strong dependency between tests, a situation commonly encountered in QTL mapping because all tests within the same chromosome will naturally show certain dependence due to genetic linkage. The permutation test is by far the most popular method used in single QTL mapping, and the score-based seems to be more suitable for multiple QTL mapping.

*Permutation*: The goal of the permutation is to build an empirical distribution of the test statistic when the null hypothesis is true. In the case of QTL mapping, the null hypothesis of no QTL effect on a trait translates into the lack of statistical association between the phenotype and genotypes of markers segregating in the population. An obvious strategy of mimicking the null hypothesis is by breaking down any relationship between phenotype and marker genotypes, as implemented in the standard permutation test of CHURCHILL and DOERGE (1994). The standard permutation test consists in shuffling the trait values, keeping the marker genotypes as in the observed data, for many times, and carrying out the genome scan on the permuted data. For each permuted data, the maximum LRT across the genome is kept. The threshold

for a given $\alpha$-level of type I error is chosen to be the $100(1-\alpha)^{th}$ percentile of the empirical distribution of the maximum LRT. This approach is more appropriate for threshold computation in single QTL model, but it is not so appropriate for threshold computation in the multiple QTL model. Two other permutation-based methods for threshold computations, the conditional empirical threshold (CET) and the residual empirical threshold (RET), that take into account the effects of QTL already in the model were also proposed by DOERGE and CHURCHILL (1996).

In the CET method, the procedure begins with the assumption of a *priori* information about a major QTL, for which the closest marker is chosen as conditioning marker, then the genotypes of the conditioning marker are used for classifying the subjects into categories. Once the categories are formed, the phenotype shuffling is done within each one of them, and CET value is computed from the empirical distribution of the maximum LRT in the permuted data. In practice, three issues may restrict using the CET. First, if the conditioning marker has many missing data, another conditioning marker should be chosen, even if it is far from the QTL, thus reducing the strength of the stratification in improving the power of the permutation test. Second, even with stratification, the markers in the same linkage group as the conditioning marker show association with the phenotype. It was suggested that the linkage group with the conditioning marker is excluded from the computations of the empirical distribution of the maximum LRT. Third, the method is computationally demanding.

In the RET method, the residuals from the model with known QTL effects are treated as a new trait, which is used for searching unknown QTL effects and for computing the threshold using the standard permutation test of CHURCHILL and DOERGE (1994). An assumption of the RET method is that the model is correct. If this assumption is violated, the results may be misleading.

*Score-based Threshold*: The standard permutation test, CET and RET seem not to be very well adapted for computing adequate threshold values to identifying multiple QTL. An alternative model-based method, attractive both in terms of computation

burden and statistical properties, for computing threshold values in multiple QTL search has been proposed by ZOU *et al.* (2004), the score-base threshold. C. Laurie, S. Wang, L. A. Carlini-Garcia and Z-B. Zeng (unpublished) extended the score statistic for testing hypothesis in the MIM method.

Under some regular conditions, the score and LRT statistics are asymptotically equivalent in large sample. But, an interesting characteristic of the score statistic is that it can be approximated by a sum of independent random numbers or vectors, depending on whether one or more parameters are tested at the same time, respectively. Motivated by this characteristic and based on the decomposition of the score function in COX and HINKLEY (1974), ZOU *et al.* (2004) derived the large-sample distribution of the score statistic for genome-wide QTL mapping.

The re-sampling algorithm for the computation of the score-based threshold consists of four steps. Before detailing the procedure, let's re-iterate the hypotheses under investigation. In genome-wide scan a putative QTL is assumed at every position $\lambda \in \zeta$ and its effect significance (main or epistatic effect) is tested against the null of no effect. Suppose a model with $m - 1$ main QTL effects and $p$ epistatic effects and we are scanning for a putative $m^{th}$ QTL. Let $l = \lambda$ denotes the testing position of the putative QTL coming into the model and $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \cdots, \lambda_{(m-1)}, l)$ be the current positions of all QTL. Let $\theta_m = \beta_m$ represents the effect for the new QTL coming into the model, and let $\boldsymbol{\theta} = (\theta_1, \theta_2, \cdots, \theta_{m-1}, \theta_m, \theta_{m+1}, \cdots, \theta_s, \mu, \sigma_e^2)'$ be the column vector of all parameters, where $\theta_b = \beta_b$ for $1 \leq b \leq m$ and $\theta_b = w_b$ for $m < b \leq s = m + p$. Let $\boldsymbol{\eta} = (\theta_1, \theta_2, \cdots, \theta_{m-1}, \theta_{m+1}, \cdots, \theta_s, \mu, \sigma_e^2)'$ be the column vector of nuisance parameters. The hypotheses of interest are $H_0 : \theta_m = 0$ versus $H_1 : \theta_m \neq 0$. Let $\tilde{\boldsymbol{\eta}}$ be the MLE of $\boldsymbol{\eta}$ under $H_0$ at testing position $l$. Let $\hat{U}_i(l)$ be equation (1.3) evaluated at testing position $l$. Similarly let $\hat{U}(l) = \sum_{i=1}^{n} \hat{U}_i(l)$ and $\hat{V}(l) = \sum_{i=1}^{n} \hat{U}_i(l)\hat{U}_i'(l)$ be evaluations of $\hat{U}$ and $\hat{V}$ at testing position $l$, respectively. Then, the steps of the re-sampling algorithm are:

1. generate $n$ independent samples $z_i$ $(i = 1, 2, \cdots, n)$ from a normal distribution

with mean zero and unit variance;

2. for each testing position $l$, obtain $\tilde{\boldsymbol{\eta}}$ and compute $\hat{U}^*(l) = \sum\limits_{i=1}^{n} \hat{U}_i(l)z_i$ and $S^*(l) = \hat{U}^{*\prime}(l)\hat{V}^{-1}(l)\hat{U}^*(l)$. Then, compute $S^* = \max\limits_{l \in \boldsymbol{\zeta}}\{S^*(l)\}$;

3. repeat *steps* 1 and 2 many times, say $N$ times (re-sampling), to obtain a sequence $(S_1^*, S_2^*, \cdots, S_N^*)$;

4. the score-based threshold for a given significance $\alpha$-level is the $100(1-\alpha)$ percentile of the ascending ordered values $(S_{(1)}^*, S_{(2)}^*, \cdots, S_{(N)}^*)$.

If $\hat{U}_i(l)$ in $\hat{U}^*(l)$ and $\hat{V}(l)$ are assumed to be fixed and $z_i$ in $\hat{U}^*(l)$ to be random, then: **I** - the conditional distribution of $\hat{U}^*(l)$ on the observed data is normal with mean zero and limiting covariance as that of $\hat{U}(l)$; **II** - from **I**, it follows that the distributions of $n^{-\frac{1}{2}}\hat{U}^*(l)$ and $n^{-\frac{1}{2}}\hat{U}(l)$ are asymptotically equivalent and; **III** - from **II**, it is possible to approximate the distribution of $S(l)$ by that of $S^*(l)$ under the null hypothesis (Zou *et al.*, 2004; Lin, 2005).

A noticeable burden in the score statistic is the necessity of analytical derivatives, especially the second order derivatives, which may be very messy. But, once the derivatives expressions have been derived, the computation of the score statistic is very straightforward and a matter of computer implementation. The score-based threshold computation is quite fast, because $\hat{U}_i(l)$ and $\hat{V}(l)$ are computed only once at each genomic position and the remaining step of generating random variables can be implemented very time efficiently. In contrast, the requirement of analyzing many permuted data sets makes the permutation test very time consuming, especially if the data set is large. It is worth mentioning that the score-based procedure has been successfully implemented and shown to have very good stability and performance in the multiple QTL search in single trait analysis as well as in multiple trait analysis. **Selecting a class of models**: Very often an investigator is looking for two classes of models, the purely additive and the additive and epistatical. The class of purely additive models is appealing because of its simplicity, which reduces tremendously the

burden of model selection and facilitates the interpretability of the fitted model. However, this simplicity reduces its spectrum of applications in practice, where research results have been shown the importance of interacting loci affecting quantitative traits (for instance, WEBER *et al.* (1999); REIFSNYDER *et al.* (2000)). The class of additive and epistatical models is gradually gaining the attention of researchers due to the possibility of fitting complex models including epistatic effects.

Mathematically, the inclusion of interaction terms into the model is trivial, though in practice, fitting a realistic model with epistasis is very challenging because of the lack of robust search procedure and good criteria that statistically maintain the false discovery rate under control. In most of the current multiple QTL methods, only the subset of identified loci with main effects are considered in the search for epistasis. The search for epistatic loci without identified main effects is an ongoing topic of research.

**Model selection criteria**: The goal of model selection is to identify a subset of regressors that is actually related to the response variable $y$. The collection of criteria for model selection based on minimizing prediction error is vast in the literature of model selection. More traditional criteria are the Mallows's $C_p$ and the adjusted $R^2$. They have been reported to be liberal, in the sense of selecting oversized models. A second generation of approaches that minimizes a criterion of the form $C(\kappa) = -2\log(\hat{L}_\kappa) + D(n, P)$, where $P$ is the total of parameters in model $\kappa$, $n$ is the sample size, and $\hat{L}_\kappa$ is the likelihood of model $\kappa$, includes the Akaike's information criterion (AIC) (AKAIKE, 1969), the Bayesian information criterion (BIC) (SCHWARZ, 1978), and the Akaike's Information Criterion corrected (AICc) (SUGIURA, 1978), with $D(n, P)$ being equal to $P\log(n)$, $2P$, and $2P + \frac{2P(P+1)}{n-P-1}$, respectively. BROMAN and SPEED (2002) proposed the use of $D(n, P) = \delta P\log(n)$, and called the criterion $BIC_\delta$. The penalty parameter $\delta$, in genome-wide QTL studies, depends on the threshold $W$ of the maximum LOD distribution, and could be defined as $\delta = \frac{2W}{\log_{10}(n)}$. Alternative computer intensive methods for model selection are the bootstrap (EFRON and TIBSHIRANI, 1993; SHAO and TU, 1995), delete-one cross-validation (CV) and

delete-$d$ cross-validation (SHAO and TU, 1995). The delete-one cross-validation is conservative in the sense that it tends to choose a model with too many regressors (SHAO and TU, 1995).

The delete-d CV approach is an alternative to overcome the conservativeness and the inconsistency of the delete-one CV. In the delete-d CV method, the entire data set is divided into two distinct subsets of data: a construction data with $n-d$ observations $(y_i^c, \boldsymbol{M}_i^c)$ and a validation data with $d$ observations $(y_i^v, \boldsymbol{M}_i^v)$, where the superscripts $c$ and $v$ stand for construction and validation, respectively. The construction data is used for estimating the model parameters and the validation data is used for assessing the prediction error of the fitted model. The consistency of the delete-d CV, under certain weak conditions, is guaranteed if and only if $\frac{d}{n} \to 1$ and $n - d \to \infty$, i.e., the size of the validation set must be much larger than the construction set (SHAO and TU, 1995).

The prediction error in the bootstrap methods is asymptotically equivalent to the delete-one CV, therefore, it inherits its inconsistentness and conservativeness (SHAO and TU, 1995). A consistent bootstrap estimator can be obtained by following the same strategy of the delete-d CV, i.e., one forms bootstrap replicates of size $k$ $(k < n)$. If $k$ is chosen such that $\frac{k}{n} \to 0$, then the new bootstrap estimator is asymptotically equivalent to the delete-d CV with $d = n - k$ (SHAO and TU, 1995).

**Searching through the space of models**: The space of models may be very large, for instance, even for the class of additive models with as few as 10 regressors, the model space would have $2^{10} = 1024$ models. Therefore, some efficient searching technique on the subset of model space must be employed in the hope of finding good models efficiently. There are many strategies in the literature ready for use in QTL mapping, among them the backward elimination, forward selection, and stepwise selection seem to appear as favorites. Depending on the size of the model space, the subset search may reduce tremendously the burden in computations. However, the saving in computations may come with the price of perhaps missing good models.

In backward elimination, in the first step, the full model is fitted. In the subsequent

steps, one at the time, the regressor with smallest regression sum of squares in the model is dropped. In forward selection, from the initial model with no regressors, one adds, one at the time, the regressor with the largest regression sum of squares into the model. In both methods, a sequence of nested models is created at the end of the algorithm's execution.

The stepwise selection is a general strategy that allows moves in either way, eliminating or adding regressors within the same step of search. The backward stepwise selection begins with the full model, and then at each step the least significant regressor is dropped and all regressors dropped, but the one dropped in the current step, are re-considered for re-entering the model. Usually, two different significance levels are used, one for deleting and other, more stringent, for adding regressor. In the forward stepwise selection, regressors once added to the model may be deleted if they are no longer significant.

In backward elimination once a regressor is dropped, it will be excluded from all subsequent models, and in the forward selection once a regressor is added, it will remains in all further models. On the other hand, in the stepwise selection, a regressor may enter or leave the model in different steps. Therefore, the number of models visited during the search may be higher in stepwise selection, which increases the chances of finding more good models (MILLER, 2002; BROMAN and SPEED, 2002).

The score-based threshold can also be used as a criterion to build and refine models with many QTL. Starting with a model with no QTL effect one can select significant putative QTL based on the score-based threshold, add them to the model, and then further refine the model, by including or excluding QTL effects. An algorithm, analogue to the algorithm described in ZENG *et al.* (1999), to build an initial model and to refine it upon using the score-based threshold criterion could be as follows: *Forward selection*– assuming that model (1.1) starts with no QTL, one QTL is added at each step of the forward selection. In the $m^{th}$ step of the forward selection, one could assume a putative QTL at every position $l \in \boldsymbol{\zeta}$, but avoiding positions within the 5 cM neighboring regions of the $m-1$ QTL in the model already identified at

previous steps of the forward selection and compute the MLE of each parameter in the model with the QTL at positions $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \cdots, \lambda_{m-1}, l)$. For each position $l$, the LRT statistic can be used for testing the significance of the effect of the $m^{th}$ putative QTL. A putative QTL at the maximum LRT statistic over all positions $l$ is added to the model if the LRT statistic is larger than the score-based threshold. *Model optimization*: in turns, the positions of all QTL identified in the forward selection are updated. One could pick a QTL in the model, and hold the other QTL in the model fixed at the positions that they were found in the forward selection. The effects of the picked QTL are then removed from the model and a new search for a QTL is done within the region delimited by its two neighboring QTL, avoiding 5 cM from the neighbor QTL (the search is performed until the end of the chromosome if no neighbor QTL is found on either side of the picked QTL). The new position of the picked QTL is set to the position of the maximum LRT statistic within the searched region and all parameters in the model are updated. The effect of the picked QTL can then be excluded from the model if it shows no significant effect. This procedure is repeated until all positions of QTL in the model are updated.

We praise for running expertise-based analysis, in which, judgements based on the investigator's expertise is launched when selecting a subset of models, that ultimately are both biologically concise and fit the data well. In practice, perhaps, one may take advantage of a *priori* information to fit an initial model, and then use stepwise selection to include or exclude QTL effects. Once a model has been built, one may test each effect individually, excluding the least significant, and finally performing the model optimization. Another alternative could be to iterate stepwise selection and model optimization.

The score- and permutation-based thresholds for the multiple QTL mapping methods of KAO *et al.* (1999) and GARCIA *et al.* (2008), as well other common methods for mapping QTL, have been implemented in the Windows QTL Cartographer computational tool for mapping QTL (WANG *et al.*, 2007). This computational tool also includes many other functions for importing and exporting a variety of data formats,

for genome-wide scan, and graphical tools for presenting and summarizing results. Windows QTL Cartographer is freely available at http://statgen.ncsu.edu/qtlcart/WQTLCart.htm. Some of the model selection criteria presented previously are also implemented in Windows QTL Cartographer within the MIM procedure.

### 1.3.2   Multiple interval mapping: Bayesian methods

In this section we review QTL inference within the Bayesian framework. For model (1.1) with $m$ QTL, let's define $\boldsymbol{G}$ the $n$ by $m$ unknown matrix of QTL genotypes (each row of $\boldsymbol{G}$ has the $m$ QTL genotypes for a given subject) and $\boldsymbol{\psi}$ the unknown parameters in the model, which is either $\boldsymbol{\psi} = (\boldsymbol{\beta}, \mu, \sigma_e^2, \boldsymbol{\lambda}, \boldsymbol{G})$ or $\boldsymbol{\psi} = (\boldsymbol{\beta}, \mu, \sigma_e^2, \boldsymbol{\lambda}, \boldsymbol{G}, m)$, according to whether $m$ is known or unknown parameter, respectively. We assume that genotypes of genetic markers ($\boldsymbol{M}$) and distance between markers ($\boldsymbol{\mathcal{R}}$) are known. Bayesian framework combines the likelihood of the data, $p(\boldsymbol{y}|\boldsymbol{\psi}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})$, with prior knowledge of the unknowns, $p(\boldsymbol{\psi})$, through an application of Bayes theorem to produce the joint posterior distribution over all unknowns, $p(\boldsymbol{\psi}|\boldsymbol{y}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})$, which is then used for inference. Inference regarding each unknown is based on its marginal distribution, which can be obtained from the joint posterior distribution by integrating over the other unknowns. Analytical integration of the joint posterior distribution may be cumbersome or even impossible because the posterior distribution lives in high-dimensional product space. When analytical close form expressions of posterior probabilities are unaffordable, a Markov Chain Monte Carlo (MCMC) technique can be used to form a Markov Chain whose stationary distribution is the posterior distributions of the unknowns. Several standard samplers used in MCMC methods are available in the literature, such as the Metropolis (METROPOLIS *et al.*, 1953), the Metropolis-Hastings (HASTINGS, 1970), the Gibbs (GEMAN and GEMAN, 1984), and the Reversible Jump (GREEN, 1995). A good sampler spends more time in the regions of higher posterior probabilities of the product space, while still visiting regions of lower probabilities. The posterior distribution of the unknowns can be written as

follows:

$$p(\boldsymbol{\psi}|\boldsymbol{y}, \boldsymbol{M}, \boldsymbol{\mathcal{R}}) = \frac{p(\boldsymbol{y}|\boldsymbol{\psi}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})p(\boldsymbol{\psi})}{\int p(\boldsymbol{y}|\boldsymbol{\psi}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})p(\boldsymbol{\psi})d\boldsymbol{\psi}}$$

where, the integral in the denominator is over as many dimensions as the number of parameters in $\boldsymbol{\psi}$.

Several Bayesian methods of QTL inference have been proposed in the literature, some treat the number of QTL in the model as a fixed quantity and others as a random variable from state-to-state in the MCMC (SORENSEN and GIANOLA, 2002). **Methods with fixed number of QTL**: SATAGOPAN *et al.* (1996) proposed a MCMC algorithm to inferring the positions and effects of multiple QTL under the assumption of known number of QTL. The authors factorized the posterior distribution in the following manner:

$$p(\boldsymbol{\beta}, \mu, \sigma_e^2, \boldsymbol{\lambda}, \boldsymbol{G}|\boldsymbol{y}, \boldsymbol{M}, \boldsymbol{\mathcal{R}}) \propto p(\boldsymbol{y}|\boldsymbol{\beta}, \mu, \sigma_e^2, \boldsymbol{G})p(\boldsymbol{G}|\boldsymbol{\lambda}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})p(\boldsymbol{\lambda})p(\boldsymbol{\beta}, \mu, \sigma_e^2) \quad (1.4)$$

where, $\propto$ stands for proportional.

If the conditional posterior distribution of an unknown is explicit and easy to sample from, the Gibbs sampler can be used. Otherwise, the Metropolis-Hastings can be applied to obtain the desirable samples. For all parameters in $\boldsymbol{\psi} = (\boldsymbol{\beta}, \mu, \sigma_e^2, \boldsymbol{\lambda}, \boldsymbol{G})$, but for the positions of QTL ($\boldsymbol{\lambda}$), there exist conjugate priors. The sequence of states in the Markov Chain, $\boldsymbol{\psi}^{(0)}, \boldsymbol{\psi}^{(1)}, \cdots, \boldsymbol{\psi}^{(N)}$, starts at any point $\boldsymbol{\psi}^{(0)}$ with positive posterior density, and at each state, under certain rules, the unknowns are updated sequentially in the order $\boldsymbol{\lambda}$, $\boldsymbol{G}$, $\mu$, $\boldsymbol{\beta}$ and $\sigma_e^2$ (SATAGOPAN *et al.*, 1996). The collected samples of the Markov Chain after a long run are assumed to be from the target posterior distribution, and these samples are used to make inference for the parameters.

SEN and CHURCHILL (2001) developed an imputation method of QTL analysis, in which multiple versions of QTL genotypes are sampled from their posterior distribution conditional on the phenotypes and marker genotypes, and these multiple sets of QTL genotypes are then used for computing approximate posterior densities of the

parameters of interest as well as the marginal probability of the data, which may be used for model selection. This procedure arises naturally because the posterior distribution of the unknowns can be divided into two independent parts, the genetic model $p(\boldsymbol{y}|\mu, \boldsymbol{\beta}, \sigma_e^2, \boldsymbol{G})p(\mu, \boldsymbol{\beta}, \sigma_e^2)$, and the linkage model, $p(\boldsymbol{G}|\boldsymbol{\lambda}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})p(\boldsymbol{\lambda})$ (see equation 1.4). The sampling is a two-step procedure: first, for a given set of QTL positions $\boldsymbol{\lambda}$, QTL genotypes $\boldsymbol{G}$ are sampled from their conditional distribution on the marker data $p(\boldsymbol{G}|\boldsymbol{\lambda}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})$; second, each genotype is weighted by the likelihood of the phenotypic data conditional on the sampled QTL genotypes $\boldsymbol{G}$, $p(\boldsymbol{y}|\boldsymbol{G})$. Although, the derivations of expressions and the computations are framed within the Bayesian theory, the method deviates from the Bayesian approach in the model selection aspect. Single and pairwise genome scans are performed and those regions with strong evidence of QTL, i.e., regions that exceed stringent threshold from the permutation test, are selected to fit multiple QTL models which are subsequently used to make model comparisons. The genome scan is restricted to low dimensional search because the number of positions to search through in the genome-wide scan grows quickly as the number of QTL increases. This limitation in the method is because the method does not provide a procedure for sampling the set of QTL positions $\boldsymbol{\lambda} = \{\lambda_1, \lambda_2, \cdots, \lambda_m\}$, therefore, in principle all possible sets of positions of size $m$ should be used in the computations. Nevertheless, the imputation framework is able to handle covariates, missing data, non-normal phenotypic data, and simple scenarios of epistasis as well as multiple trait.

**Methods with varying number of QTL**: In SATAGOPAN *et al.* (1996) the Monte Carlo (MC) simulations are conditional on the number of QTL $m$. In order to compare models of different sizes, one has to run a MC simulation for each model size and use the Bayes factor as a measure of comparison. With the advent of the reversible jump MCMC approach (GREEN, 1995), which allows sampling from the distribution of interest when the dimension of the vector of parameters is not fixed, $m$ can be regarded as random and be sampled together with all other parameters. Several reversible jump MCMC methods of QTL inference have been proposed in the literature

(WAAGEPETERSEN and SORENSEN, 2001). It seems that SATAGOPAN and YAN-
DELL (1996) were the first in the QTL literature to propose using the reversible jump
MCMC to estimate the marginal posterior probability distribution of the unknown
number of QTL and the marginal of other parameters conditional on the number of
QTL. At each state of the Markov Chain, a QTL may be included to (birth), exclude
from (death) the model, or the number of QTL remains the same, then conditional
on the new number of QTL, the other parameters are updated under certain rules.
The sampled states can be used to obtain inference regarding the parameters of in-
terest, for instance, the frequencies of the sampled values of $m$ gives an estimate of
its marginal posterior density. Inference for the other parameters are conditional on
$m$.

Yi *et al.* (2003) extended the reversible jump MCMC for inference of the number
of QTL and their effects simultaneously when some QTL may have main effects only,
some may have epistasis effects only, and others may have main and epistasis effects.
For a given model with $m$ QTL, instead of setting the unimportant effects to zero, the
indicator variables $\gamma_r$ and $\gamma_{rl}$ ($r$ and $l \in \{1, 2, \cdots, m\}$), for main effects and epistatic
effects, respectively, are included into the model. The indicator variables $\gamma_r$ and $\gamma_{rl}$
take value one or zero, according to whether their corresponding indicated effects are
included to or removed from the model, respectively. The statistical model (1.1) can
be re-written as follows:

$$y_i = \mu + \sum_{r=1}^{m} \gamma_r \beta_r x_{ir} + \sum_{r<l} \gamma_{rl} \omega_{rl} x_{ir} x_{il} + e_i$$

If the indicator variables $\gamma_r$ and $\gamma_{rl}$ are collected into a vector $\boldsymbol{\gamma}$, and one defines
a diagonal matrix $\boldsymbol{\Gamma}$ with the elements of $\boldsymbol{\gamma}$ in its diagonal, $\boldsymbol{\Gamma} = diag(\boldsymbol{\gamma})$, then the
linear model can be rewritten in matrix notation is $\boldsymbol{y} = \boldsymbol{1}\mu + \boldsymbol{X}\boldsymbol{\Gamma}\boldsymbol{\beta} + \boldsymbol{e}$ (Yi *et al.*,
2005).

The posterior distribution of the unknowns, including the number of QTL $m$ and

the vector of indicator variables $\boldsymbol{\gamma}$, can be written as:

$$p(\mu, \boldsymbol{\beta}, \sigma_e^2, m, \boldsymbol{\lambda}, \boldsymbol{\gamma}, \boldsymbol{G} | \boldsymbol{y}, \boldsymbol{M}, \boldsymbol{\mathcal{R}}) \propto p(\boldsymbol{y} | \mu, \boldsymbol{\beta}, \sigma_e^2, m, \boldsymbol{\gamma}, \boldsymbol{G}) p(\boldsymbol{G} | m, \boldsymbol{\lambda}, \boldsymbol{M}, \boldsymbol{\mathcal{R}})$$
$$\times p(\mu, \boldsymbol{\beta}, \sigma_e^2, m, \boldsymbol{\gamma}) p(\boldsymbol{\lambda})$$

The reversible jump MCMC algorithm updates, sequentially at each state of the Markov Chain, the missing marker genotypes and missing phenotypes, the parameters $(\mu, \boldsymbol{\beta}, \sigma_e^2)$, the QTL genotypes $\boldsymbol{G}$, the QTL positions $\boldsymbol{\lambda}$, the indicator variables $\boldsymbol{\gamma}$, add a new QTL with only main effect, add a new QTL with only epistasis effect, add a new QTL with both main effect and epistasis effect with some QTL already in the model, or delete an existing QTL, and also add two QTL with epistasis between themselves or delete two existing QTL. The approach represents an advance in QTL mapping in allowing the search for epistatic QTL, which may or may not show main effect. However, there are some drawbacks, such as, challenge convergence diagnoses (YI, 2004), intensive computations, low acceptance rate for including or deleting QTL, and poor exploration of the parameter space (poor mixing) (YI *et al.*, 2005).

**Methods with fixed number of QTL using the composite space representation**: Aiming to overcome the difficulties in implementing the reversible jump MCMC approach in QTL mapping, YI (2004) proposed an alternative Bayesian method based on the composite space representation of the parameter space (GODSILL, 2001). The composite space representation gets rid of the changes in the dimensionality of the vector of parameters in the reversible jump approach by, first, creating an upper bound $U$ for the number of QTL $m$, and second, adding indicator variables for all QTL effects in the model. The composite space approach, in principle, may produce more efficient proposal designs for the parameters, because parameters out of the current model may be stored and used in the proposals when a model is visited again (YI, 2004). YI *et al.* (2005) and YI *et al.* (2007) used YI (2004)'s composite space approach to develop a method for model selection to identifying pairwise epistatic effects. While YI *et al.* (2005) use the Bernoulli distribution with prior probability of success for the indicator variables of epistatic effects $\gamma_{rl}$ independently of the indicator

variables of the main effects $\gamma_r$ and $\gamma_l$, YI *et al.* (2007) constructed probabilities for the epistatic effects ($p_e$) conditional on the indicator variables of the main effects as follows:

$$p_e = P(\gamma_{rl} = 1 | \gamma_r, \gamma_l) = \begin{cases} c_0 p & \text{if} & (\gamma_r = 0, \gamma_l = 0) \\ c_1 p & \text{if} & (\gamma_r = 0, \gamma_l = 1) \text{ or } (\gamma_r = 1, \gamma_l = 0) \\ c_2 p & \text{if} & (\gamma_r = 1, \gamma_l = 1) \end{cases}$$

where, $p = P(\gamma_r = 1) = P(\gamma_l = 1)$ and $0 \leq c_0, \leq c_1, \leq c_2 \leq 1$. Besides the inclusion of epistasis between loci, the model can also includes interaction between loci and environmental variables, for instance, sex, location, and other phenotypic traits. This prior setting has two immediate consequences: first, the constant $0 \leq c_i \leq 1$, for $i = 0$, 1 and 2, implies that the main effects are more likely to be detected than the epistasis effect; second, the constants $c_i$ may play a rule of tuning parameters, in the sense that they give different weights to the epistasis effects based on the pattern of identification of main effects. Both consequences seem reasonable, since main effects are in fact easier to detect and so are epistatic effects between loci with main effects. Therefore, an appropriate tuning of the constants $c_i$ will likely increase the chances of finding epistatic effects between pairs of loci in which none or one QTL has main effect. The composite space approach is implemented in the R/qtlbim package (YANDELL *et al.*, 2007), and a technical review discussing different choices of priors as well as manners of drawing samples from the posterior distributions in the context of mapping multiple QTL is given in YI and SHRINER (2008).

**Methods with fixed number of QTL using shrinkage**: A close related procedure to the composite space is the Bayesian shrinkage. In this approach, the full set of markers is included into the linear model, and the effects of markers unrelated or trivially related to the trait phenotypes are shrunk toward zero ("removed"). WANG *et al.* (2005) proposed assigning each effect a normal prior with mean zero and effect-specific variance, and further assigning a non-informative Jeffrey's prior to each variance. The selective shrink is possible because of the special forms of the pos-

terior mean and posterior variance of each effect. Small effect will show relative small variance, which makes both its posterior mean and posterior variance approximately zero, therefore, forcing the posterior samples of the effect to be close to zero. YI and XU (2008) proposed a two-level hierarchical model, similar to WANG *et al.* (2005)'s method in the first level, in which each effect gets assigned a normal prior with mean zero and effect-specific variance. However, in the second level, they assign the variances two types of distributions with hyper-parameters, the double exponential and the inverse-$\chi^2$. The hyper-parameters in the prior for the variances are assumed to be unknown and sampled along with the other parameters in the MCMC algorithm. In this procedure the data have greater impact on the amount of shrinkage when compared with procedures that pre-specify the hyper-parameters. The use of exponential and the inverse-$\chi^2$ distributions for the variances of effects leads to the Bayesian version of LASSO (TIBSHIRANI, 1996) and the Student-t model, respectively. YI and XU (2008) observed that their method has faster convergency than methods using Jeffrey's prior, which includes no hyper-parameter.

**Posterior inference**: In Bayesian MCMC methods, inference for the underlying genetic architecture (number of QTL, QTL locations and effects) of quantitative traits relies on posterior samples from the posterior distribution with support on the parameter space. Regions in the parameter space that best support the data will be visited more frequently. This feature is of special interest in model selection for multiple QTL, because the data may support several models and the posterior distribution highlight such uncertainty (ZOU and ZENG, 2008). Specific hypotheses are tested using Bayes factor (BF), which is defined as the ratio of marginal probabilities of the data under the two hypotheses (KASS and RAFTERY, 1995). Assuming the data $\boldsymbol{y}$ arises under one of the two hypothesis $H_1$ and $H_2$ with density distributions $p(\boldsymbol{y}|H_1)$ and $p(\boldsymbol{y}|H_2)$, and letting the prior distribution and the posterior distribution of $H_i$ be $p(H_i)$ and $p(H_i|\boldsymbol{y})$, respectively (i=1, 2), then the Bayes factor is:

$$BF_{H_1 H_2} = \frac{p(\boldsymbol{y}|H_1)}{p(\boldsymbol{y}|H_2)}$$

Alternatively, from Bayes' theorem the BF can be re-written as follows:

$$\frac{p(\boldsymbol{y}|H_1)}{p(\boldsymbol{y}|H_2)} = \frac{p(H_1|\boldsymbol{y})/p(H_2|\boldsymbol{y})}{p(H_1)/p(H_2)}$$

and the Bayes factor is the ratio of the posterior odds of $H_1$ to its prior odds (Kass and Raftery, 1995). For example, the BF for testing model $m_1$ and $m_2$ is:

$$BF_{m_1 m_2} = \frac{p(\boldsymbol{y}|m_1)}{p(\boldsymbol{y}|m_2)}$$

When the number of QTL $m$ is regarded as fixed, there in no posterior samples of it to make inference regarding the size of the model. Satagopan $et\ al.$ (1996) proposed running the MCMC algorithm for each model size $N$ times and using the estimator of the marginal probability of model $m_i$ of Kass and Raftery (1995) for BF computations. The estimator of the marginal probability used by Satagopan $et\ al.$ (1996) is:

$$\hat{p}(\boldsymbol{y}|m_i) = \left( \frac{1}{N} \sum_{s=1}^{N} \frac{h(\boldsymbol{\lambda}^{(s)}, \boldsymbol{G}^{(s)}, \boldsymbol{\beta}^{(s)}, \mu^{(s)}, \sigma_e^{2(s)})}{p(\boldsymbol{y}|\boldsymbol{G}^{(s)}, \boldsymbol{\beta}^{(s)}, \mu^{(s)}, \sigma_e^{2(s)}) p(\boldsymbol{\lambda}^{(s)}, \boldsymbol{G}^{(s)}, \boldsymbol{\beta}^{(s)}, \mu^{(s)}, \sigma_e^{2(s)})} \right)^{-1}$$

where, $h(\boldsymbol{\lambda}, \boldsymbol{G}, \boldsymbol{\beta}, \mu, \sigma_e^2) = h(\boldsymbol{\lambda}) p(\boldsymbol{G}|\boldsymbol{\lambda}) p(\boldsymbol{\beta}, \mu, \sigma_e^2)$. The density $h(\boldsymbol{\lambda})$ is a normal density with support restricted to $0 \leq \lambda_1 \leq \lambda_2 \cdots \leq \lambda_m \leq L_c$, where $L_c$ is the length of the chromosome, mean $\overline{\lambda}$ and variance serving as tuning parameter. While it is clear how to apply this estimator to one chromosome only, it is not so straight for multiple chromosomes because the ordering $0 \leq \lambda_1 \leq \lambda_2 \cdots \leq \lambda_m \leq L_c$ and $\overline{\lambda}$ are meaningless in such situation.

Both BF and probability value (p-value) have to be compared to a given threshold to make the final decision regarding which hypothesis is more likely correct. However, while p-value is calculated only under the null hypothesis, BF is computed under both hypotheses (Kass and Raftery, 1995).

Consistent estimates of parameters as well as their density distributions can be ob-

tained with the posterior samples. For example, Yi *et al.* (2005) proposed estimating the main effect ($\beta$) of any QTL in $\boldsymbol{\zeta}$ or chromosome region $\Delta$ by:

$$\hat{\beta}(\Delta) = \frac{1}{N} \sum_{s=1}^{N} \sum_{r=1}^{m} I(\lambda_r^{(s)} \in \Delta, \gamma_r^{(s)} = 1) \beta_r^{(s)}$$

where, $I(x)$ is the indicator function, which takes value one if $x = x_0$ and zero otherwise.

Other estimates, such as the posterior probability of inclusion of a particular QTL $\zeta_h \in \boldsymbol{\zeta}$, $p(\zeta_h | \boldsymbol{y})$, and its Bayes factor, $BF(\zeta_h)$, are as follows:

$$\hat{p}(\zeta_h | \boldsymbol{y}) = \frac{1}{N} \sum_{b=1}^{N} \sum_{r=1}^{m} I(\lambda_r^{(s)} = \zeta_h, \gamma_r^{(s)} = 1) \quad and \quad BF(\zeta_h) = \frac{\hat{p}(\zeta_h | \boldsymbol{y})}{1 - p(\zeta_h)} \frac{1 - \hat{p}(\zeta_h | \boldsymbol{y})}{p(\zeta_h)}$$

Both $p(\zeta_h | \boldsymbol{y})$ and $BF(\zeta_h)$ can be used to profile the QTL activity across the entire genome (Yi *et al.*, 2005).

**Prior choice**: As noticed, the posterior samples which are used for inference regarding the unknown parameters combine information from the observed data and prior knowledge expressed in terms of prior probability distributions. The choice of such priors for each unknown is not so simple and may be crucial for reliable posterior inference. There exist a consensus that a *priori* information from previous related studies should anchor such choices. Experimental results have shown evidences that quantitative traits are genetically regulated by few major loci (large effects) and many minor loci (small effects) (Paterson *et al.*, 1988; Shrimpton and Robertson, 1988; Mackay, 1996). Hoeschele and VanRaden (1993) suggested that the results of Shrimpton and Robertson (1988) regarding the frequency of effects of loci affecting bristle number in *Drosophila Melanogaster* could be modeled using the exponential distribution, and thus Hoeschele and VanRaden (1993) proposed exponential prior for QTL effects. The normal distribution has been widely used as prior for the QTL effects, in part because it is conjugate under normality assumption of the data, therefore, making the derivations of posterior distributions much easier.

As important as the choice of priors for the QTL effects, is the choice of their locations throughout the genome. A common approach used regarding the QTL positions is to give all positions in the genome the same prior probability of harboring a QTL, i.e., the uniform distribution on the entire genome (SATAGOPAN *et al.*, 1996).

Numerous Bayesian approaches for QTL mapping have been proposed, and besides their advantage of combining information from the data and prior knowledge, they also can accommodate complex models. However, the choice of prior may affect the posterior distributions of the unknowns, therefore, one must be careful when picking prior distributions. Perhaps, one could carry out some sensitivity analyses of priors to better support decision regarding their choice. Other factors that might interfere in the quality of Bayesian inference are: initial values for the Markov Chain, rule to stop the chain, and mixing of the chain. Moreover, estimates of standard error of parameters from MCMC simulations are not so trivial.

## 1.4   Mapping multiple QTL on multiple trait

In genetic experiments, although in many cases multiple trait are measured in the same subject, the common rout taken for identifying QTL is to proceed with single trait analyses. However, the single trait analysis does not take advantage of the information in the data provided by the existence of phenotypic correlation between traits. The source of phenotypic correlation between traits is both genetical and environmental. In quantitative genetics, it is of practice to decompose the phenotypic variance-covariance between traits ($\boldsymbol{\Sigma}_p$) into genetic variance-covariance ($\boldsymbol{\Sigma}_g$) and residual variance-covariance ($\boldsymbol{\Sigma}_e$), in other words, $\boldsymbol{\Sigma}_p = \boldsymbol{\Sigma}_g + \boldsymbol{\Sigma}_e$ (LYNCH and WALSH, 1997). The cause of genetic correlation between traits is the existence of loci affecting the traits simultaneously (pleiotropic loci) and loci in linkage disequilibrium. Multiple measurements can also arise when one trait is measured in multiple environments. Since we can regard the expression of a trait in different environments as different trait status (FALCONER, 1952), the multivariate analysis is also very useful within

this context. If the same set of genotypes is evaluated in different environments, and significant difference of the phenotypic values is observed between environments, it is because of genes having different responses according to the environment, i.e., there is gene by environment interaction.

The literature in multivariate QTL analysis is not as extensive as the single trait one. Perhaps, this lack of research in multivariate QTL analysis is because of its complexity. Nevertheless, some approaches have been proposed. JIANG and ZENG (1995) extended the composite interval mapping (ZENG, 1993, 1994) to multiple trait analysis within the maximum likelihood framework. Besides the advantages of possibly testing for pleiotropy versus close linkage and QTL by environment interaction, the proposed method may improve, in some circumstances, the power of detecting QTL and the accuracy in estimating the positions of QTL. WELLER *et al.* (1996) and MANGIN *et al.* (1998) have both proposed the use of the method of canonical transformation of the multivariate data and the subsequent use of single trait analyses on the independent canonical variables. Although the canonical transformation technique is a powerful tool for creating independent new variables, in QTL analysis this independency may not be satisfactory because of the way that QTL affect traits. While some QTL may affect one trait only, others may affect two or more traits depending on the number of traits being analyzed jointly, i.e., the QTL may have different patterns of pleiotropy. This lack of pattern makes impossible to build independent canonical variables in which all QTL are non-pleiotropic. Creating canonical variables with all non-pleiotropic QTL is possible only if the phenotypic and genotypic correlations are equal and the pattern of pleiotropy are identical for all QTL (KNOTT and HALEY, 2000). KNOTT and HALEY (2000), following the same reasoning as in HALEY and KNOTT (1992), extended the single trait least square method to multivariate least square. A disadvantage that both methods of JIANG and ZENG (1995) and KNOTT and HALEY (2000) share is the assumption of multinormality of the traits. If this assumption is not plausible, for example when traits are count, proportion or positive continuous, these methods may lack robustness. Having this in mind, LANGE

and WHITTAKER (2001) developed a generalized estimation equations (GEE) method that deals with non-normal traits. The advantage of GEE is the possibility of estimating the parameters without a distributional assumption on the data. However, a drawback of the method is the need for supplying a variance-covariance matrix to begin with.

With the multivariate analysis some interesting questions, as the ones that follows, could be tackled: (A) Is a QTL pleiotropic or are there close linked QTL affecting multiple trait? (B) Does a QTL show genotype by environment interaction? In this section we introduce a model parametrization that allows for addressing question (A), then we describe some hypotheses setting that allows one to address question (B), and end with a discussion regarding threshold and model choice in the context of mapping multiple QTL on multiple trait.

Assuming that we have a BC population and multiple trait ($t = 1, 2, \cdots, T$) are measured in each subject of a sample of size $n$, the statistical model for multiple trait and multiple QTL analysis relates the phenotype of each subject $i$, $y_{it}$, to the coded variables $x_{ir}$ for each QTL $r$ ($r = 1, 2, \cdots, m$) affecting the traits. The explanatory variables are defined according to the Cockerham genetic model (JIANG and ZENG, 1995). A linear model with a subset $p$ of all pairwise interactions between two loci can be written as:

$$y_{ti} = \mu_t + \sum_{r=1}^{m} \beta_{tr} x_{ir} + \sum_{r<l}^{p} w_{trl} x_{ir} x_{il} + e_{ti} \tag{1.5}$$

where, for each $t$, $\mu_t$ is the mean of trait $t$ ($t = 1, 2, \cdots, T$), $\beta_{tr}$ is the effect of QTL $r$ ($r = 1, 2, \cdots, m$), $w_{trl}$ is the epistasis between loci $r$ and $l$, and $e_{ti} \sim N(0, \sigma_{e_t}^2)$. For each subject $i$, let $\boldsymbol{y}_i = (y_{1i}, y_{2i}, \cdots, y_{Ti})'$ be the $T$ by 1 vector of trait values, $\boldsymbol{X}_i$ be the $m + p$ by 1 incidence matrix, $\boldsymbol{e}_i$ be the $T$ by 1 vector of residuals, $\boldsymbol{\mu} = (\mu_1, \mu_2, \cdots, \mu_T)'$ be the $T$ by 1 vector of mean. For each $r$ and $l$, let $\boldsymbol{\beta}_r = (\beta_{1r}, \beta_{2r}, \cdots, \beta_{Tr})'$ and $\boldsymbol{w}_{rl} = (w_{1rl}, w_{2rl}, \cdots, w_{Trl})'$. We collect all the effects parameter into a $T$ by $m + p$ matrix $\boldsymbol{\mathcal{B}} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \cdots, \boldsymbol{\beta}_m, \boldsymbol{w}_1, \boldsymbol{w}_2, \cdots, \boldsymbol{w}_p)$. Then,

the linear model (1.5) in matrix form looks like:

$$\boldsymbol{y}_i = \boldsymbol{\mu} + \boldsymbol{\mathcal{B}X}_i + \boldsymbol{e}_i \tag{1.6}$$

where, $\boldsymbol{e}_i \sim MVN(\boldsymbol{0}, \boldsymbol{\Sigma}_e)$ and $\boldsymbol{\Sigma}_e$ is a positive definite symmetric residual variance-covariance matrix. We collect all $s = m + p$ vectors of effect parameters, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_e$ of model (1.6) into the vector $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}', vect(\boldsymbol{\Sigma}_e))'$, where $\boldsymbol{\theta}_b = \boldsymbol{\beta}_b'$ for $1 \le b \le m$ and $\boldsymbol{\theta}_b = \boldsymbol{w}_b'$ for $m < b \le s$, and $vect(\boldsymbol{\Sigma}_e)$ is an operator that stacks the rows of $\boldsymbol{\Sigma}_e$ into a column vector one on the top of the other and then transposes it.

Let $\boldsymbol{Z}$ be the transpose of the incidence matrix $\boldsymbol{D}$ defined previously in the single trait analysis. Analogue to the univariate model, the individual and overall likelihood of the multivariate model (1.6) with $m$ QTL are mixtures of $2^m$ multivariate normal distribution functions with different means and same variance-covariance, and mixing probabilities $p_{ij}$ ($i = 1, 2, \cdots, n$ and $j = 1, 2, \cdots, 2^m$). These mixing probabilities are defined in the same manner as previously explained in the single trait likelihood equation (1.2). The individual likelihood ($L_i$) is:

$$L_i\left(\boldsymbol{\theta} \mid \boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right) = \sum_{j=1}^{2^m} p_{ij}\phi\left(\boldsymbol{y}_i | \boldsymbol{\mu} + \boldsymbol{\mathcal{B}Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e\right)$$

where $\phi(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ is the probability density distribution of a multivariate normal random variable $\boldsymbol{z}$ with mean $\boldsymbol{\mu}_0$ and variance-covariance matrix $\boldsymbol{\Sigma}_0$. The overall likelihood ($L$) is $L\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}\right) = \prod_{i=1}^{n} L_i\left(\boldsymbol{\theta} \mid \boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right)$, where $\boldsymbol{Y}$ is a $T$ by $n$ matrix containing the traits measurements.

The expectation maximization (EM) algorithm (DEMPSTER *et al.*, 1977) solves the incomplete logarithm likelihood iteratively in terms of the unobserved complete logarithm likelihood. The E-step at the $(\nu + 1)$ iteration consists of updating the

probabilities $\pi_{ij}$ as follows:

$$\pi_{ij}^{(\nu+1)} = \frac{\phi(\boldsymbol{y}_i|\boldsymbol{\mu}^{(\nu)} + \boldsymbol{\mathcal{B}}^{(\nu)}\boldsymbol{Z}_{[\cdot,j]}, \Sigma_e^{(\nu)})p_{ij}}{\sum\limits_{j=1}^{2^m} \phi(\boldsymbol{y}_i|\boldsymbol{\mu}^{(\nu)} + \boldsymbol{\mathcal{B}}^{(\nu)}\boldsymbol{Z}_{[\cdot,j]}, \Sigma_e^{(\nu)})p_{ij}}$$

The M-step at the $(\nu + 1)$ iteration consists of updating the parameters as follows:

$$\boldsymbol{\mu}^{(\nu+1)} = \frac{1}{n}\sum_{i=1}^{n}\left(\boldsymbol{y}_i - \sum_{j=1}^{2^m}\pi_{ij}^{(\nu+1)}\boldsymbol{\mathcal{B}}^{(\nu)}\boldsymbol{Z}_{[\cdot,j]}\right)$$

$$\boldsymbol{\mathcal{B}}_{[t,b]}^{(\nu+1)} = \frac{\sum\limits_{i=1}^{n}\sum\limits_{j=1}^{2^m}\pi_{ij}^{(\nu+1)}\boldsymbol{Z}_{[b,j]}[(y_{ti} - \mu_t^{(\nu)}) - \sum\limits_{u=1}^{b-1}\boldsymbol{Z}_{[u,j]}\boldsymbol{\mathcal{B}}_{[t,u]}^{(\nu)} - \sum\limits_{u=b+1}^{s}\boldsymbol{Z}_{[u,j]}\boldsymbol{\mathcal{B}}_{[t,u]}^{(\nu)}]}{\sum\limits_{i=1}^{n}\sum\limits_{j=1}^{2^m}\pi_{ij}^{(\nu+1)}\boldsymbol{Z}_{[b,j]}^2}$$

$$\boldsymbol{\Sigma}_e^{(\nu+1)} = \frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{2^m}\pi_{ij}^{(\nu+1)}(\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu)} - \boldsymbol{\mathcal{B}}^{(\nu)}\boldsymbol{Z}_{[\cdot,j]})(\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu)} - \boldsymbol{\mathcal{B}}^{(\nu)}\boldsymbol{Z}_{[\cdot,j]})'$$

For any small positive number $\epsilon$, a stoping rule for the iterations can be defined as $\log_e L(\boldsymbol{\theta}^{(\nu+1)}|\boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}) - \log_e L(\boldsymbol{\theta}^{(\nu)}|\boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}) < \epsilon$. More efficient estimators for all parameters in model (1.6) have been derived (L.C. Silva and Z-B. Zeng, unpublished) using a expectation-conditional maximization algorithm (MENG and RUBIN, 1993) as well as a hybrid algorithm combining the expectation-conditional maximization and Newton-Raphson methods (RAI and MATTHEWS, 1993; AITKIN and AITKIN, 1996).

**Pleiotropy versus close linkage**: As previously stated, an advantage of multiple trait analysis is the possibility of testing for a single QTL affecting multiple trait versus the alternative of two or more closely linked non-pleiotropic loci (or, two or more closely linked pleiotropic loci). For instance, suppose we have measurements of two traits and a total of three non-epistatic QTL at positions $\lambda_1$, $\lambda_2$ and $\lambda_3$ in the genome. The multiple trait multiple QTL pleiotropic model for a subject $i$ would

look like:

$$
\begin{pmatrix} y_{1i} \\ y_{2i} \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} + \begin{pmatrix} \beta_{11} & \beta_{12} & \beta_{13} \\ \beta_{21} & \beta_{22} & \beta_{23} \end{pmatrix} \begin{pmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{pmatrix} + \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix} \tag{1.7}
$$

The model above assumes that all QTL have the same pattern of pleiotropy, but instead, suppose we want to test whether a QTL, say the last QTL in the model above, has effect on multiple trait against the alternative of two linked non-pleiotropic loci. The model with two pleiotropic (positions $\lambda_1$ and $\lambda_2$) and two nonpleiotropic QTL (positions $\lambda_3$ and $\lambda_4$) for a subject $i$ would look like:

$$
\begin{pmatrix} y_{1i} \\ y_{2i} \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} + \begin{pmatrix} \beta_{11} & \beta_{12} & \beta_{13} & 0 \\ \beta_{21} & \beta_{22} & 0 & \beta_{24} \end{pmatrix} \begin{pmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{pmatrix} + \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix} \tag{1.8}
$$

Or, suppose we want to test whether the last two QTL in the model (1.8) are both pleiotropic, against the alternative that they are two linked non-pleiotropic QTL. The model with four pleiotropic QTL for a subject $i$ would look like:

$$
\begin{pmatrix} y_{1i} \\ y_{2i} \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} + \begin{pmatrix} \beta_{11} & \beta_{12} & \beta_{13} & \beta_{14} \\ \beta_{21} & \beta_{22} & \beta_{23} & \beta_{24} \end{pmatrix} \begin{pmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{pmatrix} + \begin{pmatrix} e_{1i} \\ e_{2i} \end{pmatrix} \tag{1.9}
$$

Many hypotheses can be formulated and tested using model (1.5), for example, the hypotheses of model (1.7) versus (1.8) can be stated as:

$$
H_0 : \lambda_3 = \lambda_4
$$
$$
H_1 : \lambda_3 \neq \lambda_4
$$

and the hypotheses of model (1.8) versus (1.9) can be stated as:

$$H_0 : \beta_{14} = \beta_{23} = 0$$
$$H_1 : \beta_{14} \neq 0 \text{ and } \beta_{23} \neq 0$$

A couple of alternative tests have been proposed to make a statistical decision regarding model choice between a pleiotropic and a close linkage non-pleiotropic model. The LRT statistic, with a threshold from the $\chi^2_{d_1 - d_0}$ distribution was proposed by JIANG and ZENG (1995), where $d_1$ and $d_0$ are the total of parameters in the full model (under H$_1$) and reduced model (under H$_0$), respectively. A disadvantage of the LRT is the assumption of asymptotically equivalency to a $\chi^2_{d_1 - d_0}$ distribution, which may not be true since two models may not be nested within each other, as in the example previously stated regarding models (1.7) and (1.8). A parametric bootstrap strategy to empirically estimate the threshold of the distribution of the approximate LRT between the pleiotropic model and close linkage model was proposed by KNOTT and HALEY (2000). Each bootstrap replicate consists of the original marker data and trait values simulated from the pleiotropic model. This approach overcomes the asymptotic assumption for threshold computation. However, there is assumption that the pleiotropic model is correct for generating the null distribution.

A non-parametric bootstrap method that does not involve multivariate analysis was proposed by LEBRETON *et al.* (1998). It consists of drawing samples with replacement from the original data followed by single trait analyses. For each bootstrap sample, the distance between the QTL locations is computed and used to build an empirical 95%-confidence interval. If the confidence interval includes zero, the null hypothesis of pleiotropy is not rejected. This procedure has been tested with two traits only, and its application in cases where more than two traits are analyzed jointly is not so clear. KNOTT and HALEY (2000) showed in a simulation study that the non-parametric method was less powerful than the parametric, and both of them were conservative in rejecting the hypothesis of pleiotropy.

**QTL by environment interaction**: The possibility of testing for QTL by envi-

ronment interaction arises as another advantage of the multivariate analysis. There are two situations in which we are able to study the differential expression of a QTL. One, which has already been mentioned, is when the same set of genotypes are evaluated phenotypically in different environments (design I), and the other is when the phenotypic evaluations are done in different sets of genotypes, or subjects, in different environments (design II) (JIANG and ZENG, 1995). We regard the model for analysis of data in design II as multiple populations model, and thus we shall omit further discussion about it while talking about the multiple trait analysis in this review.

Let's re-iterate that in design I we regard the expression of a trait in different environments as different trait status, therefore, the index $t$ $(t = 1, 2, \cdots, T)$, which was previously defined to index traits, is regarded as the environment index in what follows. With this in mind, testing whether a QTL $r$ equally affects a trait or not in a subset $S$ $(S \in T)$ of environments, involves testing the hypotheses:

$$H_0 : \beta_{tr} = \beta_r \ \forall \ t \in S$$
$$H_1 : \beta_{tr} \neq \beta_r \ \text{for some} \ t \in S$$

The LRT statistic may be used to evaluate the hypotheses above. The critical value for the test can be obtained from the $\chi^2$ distribution with degrees of freedom being the difference in the number of parameters between the full model ($H_1$) and the reduced model ($H_0$).

**Threshold and model choice**: The usefulness of multiple trait analysis has been demonstrated in several studies (JIANG and ZENG, 1995; KNOTT and HALEY, 2000; EAVES *et al.*, 1996; WU *et al.*, 1999; LEBRETON *et al.*, 1998). However, procedures for building a model are yet scarce. As raised by KNOTT and HALEY (2000), this scarcity may be because the strategy to build a model is very likely to depend on both the underlying genetic architecture and the goals of the investigator. One may choose a bottom up approach, in which a series of single trait analyses is carried out with follow up multiple trait analysis based on the *priori* information obtained in the single trait analyses. In contrast, others may option to begin straight with the

multiple trait analysis and build a model with forward selection using some criterion for controlling the inclusion of unimportant QTL (top down approach). Or yet, one may choose to narrow down the multiple trait analysis to cluster of related variables. The cluster could be chosen based on phenotypic correlation, genotypic correlation or some other criterion, like gene pathway.

Assuming two traits and at most one QTL affecting a trait (the QTL is either pleiotropic or affects only one trait), KNOTT and HALEY (2000) studied the behavior of three procedures for identifying QTL, *single trait approach*, *linkage model approach* and *pleiotropic model approach*. In the *single trait approach*, single trait analyses were performed in each trait followed by multiple trait analysis for testing the hypothesis of one pleiotropic QTL versus two linked loci, each affecting only one trait. In the *linkage model approach*, a linkage model was fitted with a putative QTL for each trait at each position of genome (two-dimensional search) and tested against the null of no effect of both loci. Then, at the position of maximum test statistic, the effects of each non-pleiotropic QTL were tested against the null of zero effect. Finally, if the two effects were significant, the test of pleiotropy versus close linkage model was performed. As for the *pleiotropic model approach*, the pleiotropic model was fitted at each position in the genome and tested against the null of no effect. At the position of maximum test statistic, each trait effect was tested individually and if both were significant, the pleiotropic and close linked QTL model were tested. A general conclusion, as one might have expected, is that the best approach depends on the underlying genetic architecture of the traits under analysis. For example, when no residual correlation was simulated, all approaches' performances were very similar. In contrast, when the residual correlation was set to 0.75, the pleiotropic model performed the best. Also, when QTL were simulated with effects on both traits, the multiple trait model were more significant than the single trait model. MAIA (2007) chose the bottom up procedure to build an initial multiple trait multiple QTL model. Each trait was analyzed with the multiple interval mapping (KAO *et al.*, 1999), and the number and positions of the identified QTL were used as initial values in the multiple trait

multiple QTL model. The permutation test was adopted to obtain the genome-wide threshold in the single trait analyses. After sketching an initial model, the positions of the QTL were adjusted with the pleiotropic versus close linkage hypothesis testing.

Two weakness in KNOTT and HALEY (2000) study are, the assumptions of two traits and at most one QTL affecting each trait. The assumption of one QTL in the model permitted them to empirically estimate the threshold for a given genome-wide significance level in both the linkage and pleiotropic model approaches, by simulating several replicates under the null hypothesis of no QTL effects on each trait. In real data analysis, their suggestion is to use permutation to empirically estimate the threshold. MAIA (2007), as well, applied the permutation test in the multiple interval mapping. However, as mentioned previously in the review, the permutation test for multiple QTL model has some drawbacks.

While many Bayesian methods have been proposed for mapping QTL on single trait, just recently BANERJEE *et al.* (2008) seems to have proposed the first Bayesian method for mapping QTL on multiple trait. Their method allows for mapping sets of QTL specific to each trait through the concept of *seemingly unrelated regression* (ZELLNER, 1962), in which each trait is allowed to have its on set of QTL. Without any doubts the method represents an advance in mapping multiple QTL on multiple trait, however, no criterion has been proposed for assessing significance of putative QTL.

The lack of a good criterion for assessing significance of putative QTL when mapping multiple QTL on multiple trait has motivated us to extend the score statistic of ZOU *et al.* (2004) for threshold computation in the context of multiple trait (to be published elsewhere) and this technique will be included into the already available multiple trait multiple QTL mapping modulus in the Windows QTL Cartographer (WANG *et al.*, 2007).

In this paper we reviewed statistical methods for QTL mapping of non-dynamic traits from populations derived from inbred line crosses. There is also a vast literature of statistical methods for QTL mapping of dynamic (functional, longitudinal)

complex traits, for instance, MA *et al.* (2002), ZHAO *et al.* (2005), WU and LIN (2006), YANG *et al.* (2006), YANG and XU (2007), YAP *et al.* (2009) and LIU and WU (2009). Likewise, there is also a rich literature of statistical methods for mapping QTL from data of outbred populations (for instance, half-sibs and full-sibs), family-based populations (pedigree data) and observational populations (association mapping) (see, for instance, SLATE (2005), BALDING *et al.* (2007), NEALE *et al.* (2007), and MCCARTHY *et al.* (2008)).

## 1.5  Final considerations

In this paper we reviewed the multiple interval mapping both within the maximum likelihood and Bayesian frameworks. The main advantage throughout the development of multiple interval mapping is the improvement in identification of putative loci due to fitting of more complex models that ultimately use the information available in the data more efficiently through information of many marker intervals simultaneously. Another advantage brought by modeling multiple loci simultaneously is the possibility of evaluating epistasis between loci.

In the maximum likelihood method, the data is fitted with a mixture of distribution functions (we reviewed only the mixture of normal distributions), which requires somewhat more sophisticated approaches of parameter estimation, such as the EM algorithm (DEMPSTER *et al.*, 1977). To build a set of models with multiple loci within the maximum likelihood framework the genome is usually partitioned in small bins (or grids), say 1-cM, and the model is fitted at every position in the genome and the LRT is used for the assessment of significance of effects at each position. The position within the genome with maximum LRT is considered as harboring a putative QTL if the LRT at the position exceeds a predefined criterion. A good criterion must correct for the multiplicity of testing and take into account the correlation between tests due to the genetic linkage. The permutation test (CHURCHILL and DOERGE, 1994; DOERGE and CHURCHILL, 1996) is by far the most widely used method for

threshold estimation in single trait single QTL analysis. However, it might lack robustness when models with multiple loci are fitted. The score-based threshold (ZOU *et al.*, 2004) has been empirically shown to be an alternative to the permutation test for computing the threshold in single trait multiple interval mapping (C. Laurie, S. Wang, L. A. Carlini-Garcia and Z-B. Zeng, unpublished).

The Bayesian paradigm provides us with tools for building a set of complex models. However, prior distributions need to be chosen, and their choice might have impact in the posterior inference. Moreover, the sampling of parameters in high-dimensional space and the diagnosis of convergence have been reported to be challenging, especially when the number of loci in the model is treated as an unknown parameter and it is sampled along with the other model parameters using the reversible jump MCMC algorithm of GREEN (1995) (YI, 2004; YI *et al.*, 2005).

In genetic experiments, although in many cases multiple trait are measured in the same subject, analyses of single traits have been the main stream for the purpose of QTL identification. However, the single trait analyses do not take advantage of the information in the data regarding the existence of genetic and environmental correlation between traits. Within the multiple trait analysis framework we reviewed interesting hypotheses and statistical methods that might allow an investigator to assess the pattern of action of loci on multiple trait, such as, testing the hypothesis of the existence of a pleiotropic QTL versus the hypothesis of close linked QTL affecting multiple trait, and testing the hypothesis that a QTL shows genotype by environment interaction. However, preceding the use of these interesting features provided by the multiple trait analysis, one may need to build a model or a set of models. KNOTT and HALEY (2000) and BANERJEE *et al.* (2008) have developed a least squares and a Bayesian method for mapping QTL multiple trait, respectively. The latter method lacks a criterion for assessing the significance level of QTL and in the former method, it was suggested the use of permutation to empirically estimate the genome-wide threshold to assess the significance level of loci effects. However, the permutation has drawbacks that might limit its robustness in multiple QTL model. Because the

multiple trait analysis is an important tool that allows an investigator to better extract information from multivariate data, therefore, revealing with more details the genetic architecture of complex traits, we have invested efforts in developing a score-based threshold for assessing significance level of QTL on multiple trait.

## 1.6 Plan of dissertation

In this dissertation we summarize our own research results on mapping QTL from inbred line crosses. We present this research in the form of four chapters. In Chapter 1, we reviewed the current status of research on statistical methods for mapping multiple QTL in single and multiple complex traits within the maximum likelihood and Bayesian frameworks. In Chapter 2, we propose a statistical model for multiple trait multiple interval mapping (MTMIM) of QTL from inbred line crosses. We also derive a set of equations for parameter estimations. We extended the score-based threshold for genome-wide evaluation of significance level of effects of QTL in the MTMIM model. We propose a strategy of forward selection using the score-based threshold as a criterion of variable selection in the MTMIM model. In Chapter 3, we implement and evaluate our MTMIM model and score-based threshold with several simulated data sets as well as with data from an experiment with *Drosophila*. In Chapter 4, we derive analytical formulae for prediction of length of confidence interval for position of QTL and for prediction of shape of the LRT around the position of QTL in multiple trait analysis in linkage maps with distinct saturation of markers.

# 2

# Multiple trait multiple interval mapping of quantitative trait loci from inbred line crosses

## 2.1 Introduction

Although in many cases multiple trait are measured in the same subject, single trait analyses have been the main stream for the purpose of QTL identification. However, single trait analyses do not take advantage of the information in the data regarding existence of genotypic and environmental correlation between traits. In Chapter 1 we reviewed interesting hypotheses and statistical methods in multiple trait analysis that might allow an investigator to assess the pattern of action of QTL on multiple trait, such as, testing the hypothesis of the existence of a pleiotropic QTL versus the hypothesis of close linked QTL affecting multiple trait, and testing the hypothesis of QTL by environment interaction. However, preceding the use of these interesting features provided by the multiple trait analysis, one may need to build a model or a set of models. KNOTT and HALEY (2000) and BANERJEE *et al.* (2008) have developed a least squares and a Bayesian method for mapping multiple

QTL in multiple trait, respectively. Both methods lack criteria for assessing the significance level of effects of QTL. In the former method, permutation test was suggested to empirically estimate the genome-wide threshold to assess significance level of effects of QTL. However, the permutation has drawbacks that might limit its robustness in multiple QTL models. Because multiple trait analysis is an important tool that allows an investigator to better extract information from multivariate data, therefore, revealing with more details the genetic architecture of complex traits, we have invested efforts in developing a score-based threshold for assessing significance level of QTL on multiple trait.

In this chapter, we propose a statistical method for multiple trait multiple interval mapping (MTMIM) of QTL from inbred line crosses. In what follows, we describe the MTMIM statistical model (Section 2.2), build the likelihood function (Section 2.3), derive parameter estimators (Section 2.4), extend the score-based threshold method of ZOU *et al.* (2004) to the MTMIM model (Section 2.5), propose a forward selection to build a model with multiple QTL using the score-based threshold as the criterion to assess the significance level of effects of QTL, and propose a model optimization procedure (Section 2.6). In the last section, we describe some alternative criteria for testing pleiotropy versus close linkage (Section 2.7).

## 2.2  Statistical model

Our statistical model for multiple trait multiple QTL inference on BC population is a linear model, in which the value of trait $t$ ($t = 1, 2, \cdots, T$), $y_{ti}$, for each $i^{th}$ subject ($i = 1, 2, \cdots, n$), is regressed on variables $x_{ir}$ ($r = 1, 2, \cdots, m$). These variables are defined according to the Cockerham genetic model (KAO and ZENG, 2002; ZENG *et al.*, 2005). For each subject $i$, $x_{ir}$ takes either value $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether QTL $r$ has genotype $QQ$ or $Qq$, respectively. The coefficient of $x_{ir}$, $\beta_{tr}$, is called the main effect of $r^{th}$ QTL on trait $t$. The linear model also includes an intercept $\mu_t$, a subset $p$ of epistatic effects ($w_{trl}$) between all pairwise interactions between

QTL ($r$ and $l \in \{1, 2, \cdots, m\}$), and a residue $e_{ti}$. The residues are assumed to be independent and identically distributed according to a normal distribution with mean zero and variance $\sigma^2_{e_t}$. The linear model is then:

$$y_{ti} = \mu_t + \sum_{r=1}^{m} \beta_{tr} x_{ir} + \sum_{r<l}^{p} w_{trl} x_{ir} x_{il} + e_{ti} \qquad (2.1)$$

For each subject $i$, let $\boldsymbol{y}_i = (y_{1i}, y_{2i}, \cdots, y_{Ti})'$ be the $T$ by 1 vector of trait values, $\boldsymbol{X}_i$ be the $m + p$ by 1 incidence matrix, $\boldsymbol{e}_i$ be the $T$ by 1 vector of residuals, $\boldsymbol{\mu} = (\mu_1, \mu_2, \cdots, \mu_T)'$ be the $T$ by 1 vector of mean. For each $r$ and $l$, let $\boldsymbol{\beta}_r = (\beta_{1r}, \beta_{2r}, \cdots, \beta_{Tr})'$ and $\boldsymbol{w}_{rl} = (w_{1rl}, w_{2rl}, \cdots, w_{Trl})'$ be column vectors of main and epistatic effects, respectively. We collect all the effect parameters into a $T$ by $m + p$ matrix $\boldsymbol{\mathcal{B}} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \cdots, \boldsymbol{\beta}_m, \boldsymbol{w}_1, \boldsymbol{w}_2, \cdots, \boldsymbol{w}_p)$, and rewrite (2.1), for a subject $i$, in matrix form as:

$$\boldsymbol{y}_i = \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{X}_i + \boldsymbol{e}_i$$

where, $\boldsymbol{e}_i$ is a random vector of length $T$ assumed to be independent and identically distributed according to a multivariate normal distribution with mean vector zero and positive definite symmetric variance-covariance matrix $\boldsymbol{\Sigma}_e$ ($MVN_T(\boldsymbol{0}, \boldsymbol{\Sigma}_e)$).

We collect all $s = m + p$ effect parameters ($m$ main and $p$ epistatic effect vectors), $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_e$ into a column vector $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}', vect(\boldsymbol{\Sigma}_e))'$, where $\boldsymbol{\theta}_b = \boldsymbol{\beta}_b'$ for $1 \le b \le m$ and $\boldsymbol{\theta}_b = \boldsymbol{w}_b'$ for $m < b \le s$, and $vect(\boldsymbol{\Sigma}_e)$ is an operator that stacks the rows of $\boldsymbol{\Sigma}_e$ into a column vector one on the top of the other and then transposes it. Motivated by the fact that a QTL may not have significant effect on all traits under analysis, we allow for the insignificant effect parameters in each vector $\boldsymbol{\theta}_b$ to be constrained to zero. Therefore, our MTMIM model allows each trait to have its own set of effect parameters, as in the *seemingly unrelated regression* model of ZELLNER (1962).

## 2.3   Likelihood function

Let's re-iterate that in order to search the entire genome for significant effects of QTL, the genome is partitioned into $H$ loci, usually at 1-cM grid. This partition is denoted by $\boldsymbol{\zeta}$. The set of positions of $m$ putative QTL, $\boldsymbol{\lambda} = \{\lambda_1, \lambda_2, \cdots, \lambda_m\}$, is assumed to be a subset of $\boldsymbol{\zeta}$ (YI *et al.*, 2005). For any subject $i$, let $\boldsymbol{M}_{[i,\cdot]}$ be the genotypic information of markers flanking the $m$ QTL, and $M_{i,L}^r$ and $M_{i,R}^r$ be the flanking markers on left and right of QTL $r$, respectively. Then, assuming no cross-over interference between marker intervals and no more than one QTL existing within a marker interval, the probability of any genotype of the form $G_j = Q_1 Q_2 \cdots Q_m$, where $Q_r \in \{QQ, Qq\}$, $r = 1, 2, \cdots, m$ and $j = 1, 2, \cdots, 2^m$, conditional on the genotypes of markers flanking the $m$ QTL is $p_{ij} = P(G_j | \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\mathcal{R}}, \boldsymbol{\lambda}) = \prod_{r=1}^{m} P\left(Q_r \mid M_{i,L}^r, M_{i,R}^r, \lambda_r\right)$, where the probabilities on the right hand side of this equation can be estimated as in JIANG and ZENG (1997) and KAO and ZENG (1997). In Table 1.2, we show how to estimate these probabilities for a BC population. The recombination frequency between QTL $Q_r$ and its left marker $(r_{LQ_r})$ has a one-to-one correspondence with the position $\lambda_r$. Conditional probabilities of two QTL lying within a single interval is shown in Table A.2 of Appendix A.

Let's also remind that $\boldsymbol{Z}$ is the transpose of the incidence matrix $\boldsymbol{D}$ defined in single trait analysis. Analogue to the univariate model, the individual and overall likelihoods of the MTMIM model with $m$ QTL are mixtures of $2^m$ multivariate normal distribution functions with different means and same variance-covariance, and mixing probabilities $p_{ij}$ $(j = 1, 2, \cdots, 2^m)$. These mixing probabilities are defined in the same manner as previously explained in the single trait likelihood equation (1.2). The individual $(L_i)$ and overall likelihoods $(L)$ are:

$$L_i \left(\boldsymbol{\theta} \mid \boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right) = \sum_{j=1}^{2^m} p_{ij} \left(2\pi\right)^{-\frac{T}{2}} |\boldsymbol{\Sigma}_e|^{-\frac{1}{2}} e^{-\frac{1}{2}\left(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]}\right)' \boldsymbol{\Sigma}_e^{-1} \left(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]}\right)}$$

$$L\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}\right) = \prod_{i=1}^{n} L_i\left(\boldsymbol{\theta} \mid \boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right) \tag{2.2}$$

where, $\boldsymbol{Y}$ is a $T$ by $n$ matrix of phenotypic trait data. In what follows, $(\ell_i)$ and $(\ell)$ represent the natural logarithm of the individual and overall likelihoods, respectively.

## 2.4   Parameter estimation

Estimation of parameters in the likelihood function (2.2) is cumbersome because mixture of distributions. The expectation-maximization (EM) (DEMPSTER *et al.*, 1977) algorithm is very popular for parameter estimation in mixture models. The EM algorithm is very simple to program, given that efficient estimators are available for the "complete-data". Moreover, the EM algorithm guarantees that the likelihood function is non-decreasing in every iteration. However, EM may show slow convergence rate if there are many missing data, and EM does not provide standard errors of parameter estimates.

Many modifications of the EM algorithm and many hybrids of EM and Gauss-Newton (GN) methods have been proposed in the literature (RAI and MATTHEWS, 1993; AITKIN and AITKIN, 1996; MCLACHLAN and KRISHNAN, 1996). GN methods are not guaranteed to converge when the logarithm likelihood is not concave, but if there is convergence its rate is usually quadratic, as opposite to the linear rate of EM. Therefore, speed of convergence of GN may be much faster than EM. In this section, we describe four algorithms to estimate parameters in the MTMIM model: EM, expectation-conditional maximization (ECM), Newtow-Raphson (NR), and a hybrid of EM and NR called generalized EM-NR (GEM-NR).

**Expectation maximization algorithm**

Let $\boldsymbol{z}_i^* = (z_{i1}^*, z_{i2}^*, \cdots, z_{i2^m}^*)'$ be a vector with information about "missing" QTL genotypes for subject $i$. Each $z_{ij}^* = 1$ if $i^{th}$ subject has genotype $G_j$, otherwise $z_{ij}^* = 0$. Let $\boldsymbol{z}^* = (\boldsymbol{z}_1^*, \boldsymbol{z}_2^*, \cdots, \boldsymbol{z}_n^*)$ be a matrix containing missing information from

all subjects. The joint distribution of observed and missing data $(\boldsymbol{y}_i, \boldsymbol{z}_i^*)$ for subject $i$ is:

$$p(\boldsymbol{y}_i, \boldsymbol{z}_i^*) = \prod_{j=1}^{2^m} \left[ \phi(\boldsymbol{y}_i | \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}, \Sigma_e) p_{ij} \right]^{z_{ij}^*}$$

where, $\phi(\boldsymbol{y}_i | \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}, \Sigma_e)$ is the probability density distribution of a multivariate normal random variable $\boldsymbol{y}_i$ with mean $\boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}$ and variance-covariance $\Sigma_e$. The mixing probabilities $p_{ij}$ are defined as previously, $p_{ij} = P(G_j | \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\mathcal{R}}, \boldsymbol{\lambda})$. The joint distribution of observed and missing data allow us to obtain the complete-data logarithm likelihood ($\ell_c$):

$$\ell_c(\boldsymbol{\theta} | \boldsymbol{Y}, \boldsymbol{z}^*) = \sum_{i=1}^{n} \sum_{j=1}^{2^m} z_{ij}^* (\log p_{ij} + \log(\phi(\boldsymbol{y}_i | \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}, \Sigma_e)))$$

The EM algorithm (DEMPSTER *et al.*, 1977) solves the incomplete logarithm likelihood iteratively in terms of the unobserved complete-data logarithm likelihood. The E-step requires computation of the expectation of the complete-data logarithm likelihood, conditional on the observed data $\boldsymbol{y}$ and evaluated at a current value of $\boldsymbol{\theta}$ (denoted here as $\boldsymbol{\theta}^{(\nu)}$) (MCLACHLAN and KRISHNAN, 1996):

$$\begin{aligned} Q_c(\boldsymbol{\theta} | \boldsymbol{\theta}^{(\nu)}) &= E_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(\nu)}} \left[ \ell_c(\boldsymbol{\theta} | \boldsymbol{y}, \boldsymbol{z}^*) | \boldsymbol{y} \right] \\ &= \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} (\log p_{ij} + \log(\phi(\boldsymbol{y}_i | \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}, \Sigma_e))) \end{aligned}$$

where,

$$\begin{aligned} \pi_{ij}^{(\nu)} &= E_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(\nu)}} \left[ z_{ij}^* | \boldsymbol{y}_i \right] \\ &= \frac{p_{ij} \phi(\boldsymbol{y}_i | \boldsymbol{\mu}^{(\nu)} + \boldsymbol{\mathcal{B}}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]}, \Sigma_e^{(\nu)})}{\sum_{j=1}^{2^m} p_{ij} \phi(\boldsymbol{y}_i | \boldsymbol{\mu}^{(\nu)} + \boldsymbol{\mathcal{B}}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]}, \Sigma_e^{(\nu)})} \end{aligned}$$

The M-step consists of maximizing the expected complete logarithm likelihood

$(Q_c)$ with respect to the unknown parameters. By taking derivatives of $Q_c$ with respect the unknown parameters (see Appendix B) and setting the first order derivatives equal to zero we can obtain the maximum likelihood estimators of all unknown parameters. In what follows, we show closed form estimators of all parameters in the MTMIM model assuming that all QTL have effects on all traits (no effect parameter is constrained to zero):

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^{n} \left( \boldsymbol{y}_i - \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{B} \boldsymbol{Z}_{[\cdot,j]} \right)$$

$$\boldsymbol{B} = \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} (\boldsymbol{y}_i - \boldsymbol{\mu})(\boldsymbol{Z}_{[\cdot,j]})' \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]} (\boldsymbol{Z}_{[\cdot,j]})' \right)^{-1}$$

$$\boldsymbol{\Sigma}_e = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} (\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{B} \boldsymbol{Z}_{[\cdot,j]})(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{B} \boldsymbol{Z}_{[\cdot,j]})'$$

The E- and M-steps are computed iteratively for many times until convergence of the likelihood function. The E-step at the $(\nu + 1)$ iteration consists of updating the probabilities $\pi_{ij}$ as follows:

$$\pi_{ij}^{(\nu+1)} = \frac{\phi(\boldsymbol{y}_i | \boldsymbol{\mu}^{(\nu)} + \boldsymbol{B}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e^{(\nu)}) p_{ij}}{\sum\limits_{j=1}^{2^m} \phi(\boldsymbol{y}_i | \boldsymbol{\mu}^{(\nu)} + \boldsymbol{B}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e^{(\nu)}) p_{ij}}$$

The M-step at the $(\nu + 1)$ iteration consists of updating the parameters as follows:

$$\boldsymbol{\mu}^{(\nu+1)} = \frac{1}{n} \sum_{i=1}^{n} \left( \boldsymbol{y}_i - \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} \boldsymbol{B}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]} \right)$$

$$\boldsymbol{B}^{(\nu+1)} = \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} (\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu)})(\boldsymbol{Z}_{[\cdot,j]})' \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} \boldsymbol{Z}_{[\cdot,j]} (\boldsymbol{Z}_{[\cdot,j]})' \right)^{-1}$$

$$\boldsymbol{\Sigma}_e^{(\nu+1)} = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} (\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu)} - \boldsymbol{B}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]})(\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu)} - \boldsymbol{B}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]})'$$

For any small positive real number $\epsilon$, a stoping rule for convergence of the likelihood can be defined as $\ell(\boldsymbol{\theta}^{(\nu+1)}|\boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}) - \ell(\boldsymbol{\theta}^{(\nu)}|\boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}) < \epsilon$.

**Expectation-Conditional maximization algorithm**

Because the assumption of unconstrained effect parameters the expressions for parameter estimation are very simple and easy to implement in the M-step described previously. This is perhaps the most attractive feature of the EM algorithm, which is possible because our "complete-data" logarithm likelihood is rather simple. If the complete-data logarithm likelihood is messy and the M-step is complex, then the EM algorithm is no longer attractive. For such cases of complicate M-step, Meng and Rubin (1993) proposed a class of generalized EM algorithm, called expectation-conditional maximization (ECM). The ECM enjoys the convergence properties of the EM while simplifying the estimation of parameters. In the ECM, a complex M-step is broken down into many simpler CM-steps, each one of them maximizes the expected complete-data logarithm likelihood conditional on some function of the parameters. Besides simplifying the M-step, the CM-step is often simpler, faster and more stable than the M-step because the conditional maximization are over spaces of smaller dimensions (Meng and Rubin, 1993). When some effect parameters in MTMIM model (2.1) are constrained to zero it is easier to implement the ECM algorithm. For instance, when estimating parameters in the model with closely linked non-pleiotropic QTL, model (1.8). Now, we describe an ECM algorithm feasible for parameter estimation in MTMIM model (2.1) when all effect parameters are unconstrained.

In the CM-step, we split the parameters into the groups $\boldsymbol{\mathcal{B}}_{[\cdot,1]}, \boldsymbol{\mathcal{B}}_{[\cdot,2]}, \cdots, \boldsymbol{\mathcal{B}}_{[\cdot,s]}$, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_e$. Parameters within the same group are estimated simultaneously, while parameters in distinct groups are estimated consecutively. The parameter estimators

can be shown to be:

$$\boldsymbol{\mu}^{(\nu+1)} = \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{y}_i - \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} \boldsymbol{\mathcal{B}}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]})$$

$$\boldsymbol{\Sigma}_e^{(\nu+1)} = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} (\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu+1)} - \boldsymbol{\mathcal{B}}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]})(\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu+1)} - \boldsymbol{\mathcal{B}}^{(\nu)} \boldsymbol{Z}_{[\cdot,j]})'$$

$$\boldsymbol{\mathcal{B}}_{[\cdot,b]}^{(\nu+1)} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} (\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu+1)} - \sum_{u=1}^{b-1} \boldsymbol{\mathcal{B}}_{[\cdot,u]}^{(\nu+1)} \boldsymbol{Z}_{[u,j]} - \sum_{u=b+1}^{s} \boldsymbol{\mathcal{B}}_{[\cdot,u]}^{(\nu)} \boldsymbol{Z}_{[u,j]}) \boldsymbol{Z}_{[b,j]}}{\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} \boldsymbol{Z}_{[b,j]}^2}$$

for $b \in \{1, 2, \cdots, s\}$.

The estimator of $\boldsymbol{\mathcal{B}}_{[\cdot,b]}$ shown previously is not appropriate if some parameters in $\boldsymbol{\mathcal{B}}_{[\cdot,b]}$ are constrained to zero. For instance, when estimating parameters in the model with closely linked non-pleiotropic QTL, model (1.8). When there exist zero-constrained effect parameters in the MTMIM model, our strategy is to update each element in $\boldsymbol{\mathcal{B}}_{[\cdot,b]}$ one at the time. Given the current estimate $\boldsymbol{\mathcal{B}}_{[\cdot,b]}^{(\nu)}$, where $\boldsymbol{\mathcal{B}}_{[t,b]}^{(\nu)} = 0$ (constrained to zero) for some $t \in \{1, 2, \cdots, T\}$, the updating equation for the unconstrained effect parameter $\boldsymbol{\mathcal{B}}_{[t,b]}$ is:

$$\boldsymbol{\mathcal{B}}_{[t,b]}^{(\nu+1)} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} \boldsymbol{\Sigma}_{e_{[t,\cdot]}}^{-1(\nu)} [(\boldsymbol{y}_i - \boldsymbol{\mu}^{(\nu)}) - \sum_{u=1}^{b-1} \boldsymbol{\mathcal{B}}_{[\cdot,u]}^{(\nu+1)} \boldsymbol{Z}_{[u,j]} - \sum_{u=b+1}^{s} \boldsymbol{\mathcal{B}}_{[\cdot,u]}^{(\nu)} \boldsymbol{Z}_{[u,j]}] \boldsymbol{Z}_{[b,j]}}{\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu+1)} \boldsymbol{\Sigma}_{e_{[t,t]}}^{-1(\nu)} \boldsymbol{Z}_{[b,j]}^2}$$

Our choice of initial values for $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_e$ are the sample mean and the sample variance-covariance, respectively, and all parameters in $\boldsymbol{\mathcal{B}}$ are initiated as zero.

It is worth mentioning that for many combinations of $i$ and $j$, the probabilities $p_{ij}$ are zero or very close to zero. Therefore, one may take advantage of sparse matrix theory to save on computation time.

**Newton-Raphson method**

The Newton-Raphson (NR) method for solving the equations (2.3) consists of approximating the gradient vector (left hand side of 2.3) by a linear Taylor's series expansion around the current fit $\boldsymbol{\theta}^{(\nu)}$ (MCLACHLAN and KRISHNAN, 1996).

$$\frac{\partial \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}} = 0 \tag{2.3}$$

The first order Taylor's approximation of the gradient vector is:

$$\frac{\partial \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}} \simeq \left.\frac{\partial \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}^{(\nu)}} + \left.\frac{\partial^2 \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^{(\nu)}} (\boldsymbol{\theta} - \boldsymbol{\theta}^{(\nu)}) \tag{2.4}$$

Equating the left hand side of (2.4) to zero and solving for $\boldsymbol{\theta}$, we obtain the updating formula for the parameters:

$$\boldsymbol{\theta}^{(\nu+1)} = \boldsymbol{\theta}^{(\nu)} + \left( -\left.\frac{\partial^2 \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^{(\nu)}} \right)^{-1} \left.\frac{\partial \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}^{(\nu)}} \tag{2.5}$$

Besides a major advantage of NR method in terms of convergence rate (when it does converge), the NR method also provides an estimate of the variance-covariance matrix of parameters in the MTMIM model at the limiting value of $\boldsymbol{\theta}$, $\boldsymbol{\theta}^*$. The inverse of the observed Fisher's information matrix provides an estimate of the variance-covariance of parameters:

$$I^{-1}(\boldsymbol{\theta}^*|\boldsymbol{Y}) = \left( -\left.\frac{\partial^2 \ell(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^*} \right)^{-1}$$

The NR method requires accurate initial values of parameters, in certain problems, in order for right convergency of the likelihood function. Moreover, the NR method has almost equally chances to move either in the direction of saddle points, local minima or local maxima (MCLACHLAN and KRISHNAN, 1996).

**Generalized EM algorithm based on Newton-Raphson methods**
The Generalized EM-Newton-Raphson (GEM-NR) methods combine the EM al-

gorithm with the NR method for maximizing the complete-data logarithm likelihood (RAI and MATTHEWS, 1993; AITKIN and AITKIN, 1996). The hybrid methods take advantage of the EM algorithm for generating an accurate starting point for the NR, and then explores the convergency rate of NR method. By introducing a step-size $\kappa^{(\nu)}$ ($0 < \kappa^{(\nu)} \leq 1$) in equation (2.5) and by having the incomplete-data logarithm likelihood replaced by the expected complete-data logarithm likelihood, we obtain a modified version for the updating equation (MCLACHLAN and KRISHNAN, 1996):

$$\boldsymbol{\theta}^{(\nu+1)} = \boldsymbol{\theta}^{(\nu)} + \kappa^{(\nu)} \left( -\left.\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^{(\nu)}} \right)^{-1} \left.\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}^{(\nu)}} \tag{2.6}$$

The advantage of using the modified version of the updating equation is that an appropriate choice of $\kappa^{(\nu)}$ guarantees the logarithm likelihood increases at each iteration. The negative of the matrix of second order derivatives in (2.6) is positive definite under usual conditions. Therefore, it has the Cholesky decomposition (2.7), where $\boldsymbol{C}$ is an upper triangular matrix.

$$\left( -\left.\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^{(\nu)}} \right)^{-1} = \boldsymbol{C}'\boldsymbol{C} \tag{2.7}$$

Let $\boldsymbol{\theta}^{(\xi)}$ be a point in the line segment from $\boldsymbol{\theta}^{(\nu)}$ to $\boldsymbol{\theta}^{(\nu+1)}$, the Taylor's expansion of the complete-data logarithm likelihood function around $\boldsymbol{\theta}^{(\nu)}$ is:

$$\begin{aligned}
Q_c(\boldsymbol{\theta}^{(\nu+1)}|\boldsymbol{Y}) - Q_c(\boldsymbol{\theta}^{(\nu)}|\boldsymbol{Y}) =& (\boldsymbol{\theta}^{(\nu+1)} - \boldsymbol{\theta}^{(\nu)})' \left.\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}^{(\nu)}} \\
&+ \frac{1}{2}(\boldsymbol{\theta}^{(\nu+1)} - \boldsymbol{\theta}^{(\nu)})' \left.\frac{\partial^2 Q_c(\boldsymbol{\theta}|y)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^{(\xi)}} (\boldsymbol{\theta}^{(\nu+1)} - \boldsymbol{\theta}^{(\nu)})
\end{aligned} \tag{2.8}$$

Plugging $\boldsymbol{\theta}^{(\nu)}$ from (2.6) into (2.8), and upon making some algebra using (2.7), we

obtain:

$$Q_c(\boldsymbol{\theta}^{(\nu+1)}|\boldsymbol{Y}) - Q_c(\boldsymbol{\theta}^{(\nu)}|\boldsymbol{Y}) = \kappa^{(\nu)} \left( \left.\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}^{(\nu)}} \right)' \boldsymbol{C}'\boldsymbol{C}\boldsymbol{B} \left.\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}^{(\nu)}} \quad (2.9)$$

where,

$$\boldsymbol{B} = \left( \boldsymbol{I} + \frac{1}{2}\kappa^{(\nu)} \left.\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{Y})}{\partial \boldsymbol{\theta}\partial \boldsymbol{\theta}'}\right|_{\boldsymbol{\theta}^{(\xi)}} \boldsymbol{C}'\boldsymbol{C} \right) \quad (2.10)$$

and $\boldsymbol{I}$ is an identity matrix.

From (2.9), we can see that as long as $\kappa^{(\nu)}$ is chosen to make (2.10) positive definite, the logarithm likelihood is guaranteed to increase at every iteration.

To guarantee that the logarithm likelihood is non-decreasing, AITKIN and AITKIN (1996) proposed to start the EM algorithm with five iterations to quickly approach the MLE and then to switch to NR until either convergence or decrease of the logarithm likelihood. If the logarithm likelihood decreases, they suggested halving the step size $\kappa$ up to five times, and if the logarithm likelihood still decreases, to return to the EM and to run five iterations and then to switch back to NR. AITKIN and AITKIN (1996) argued that their choice of running the EM algorithm for five iterations is based on previous experiences of REDNER and WALKER (1984) that 95% of the change in the initial value of logarithm likelihood until its maximum value often happens in five EM iterations.

As $\boldsymbol{\theta}^{(\xi)}$ lies in the line segment from $\boldsymbol{\theta}^{(\nu)}$ to $\boldsymbol{\theta}^{(\nu+1)}$, and $\boldsymbol{\theta}$ lives in high-dimensional space, the choice of $\kappa^{(\nu)}$ to make (2.10) positive definite may not be easy. Xue-Jun Qin and Zhao-Bang Zeng (unpublished) proposed an iterative procedure to make sure that $\kappa^{(\nu)}$ satisfies condition (2.10):

1. Let $\boldsymbol{\theta}^{(\nu)}$ be the parameter estimate in the $\nu^{th}$ iteration;

2. Set $\kappa^{(\nu)} = 1$;

3. Estimate $\boldsymbol{\theta}^{(\nu+1)}$ using (2.6) with the first and second order derivatives of $Q_c(\boldsymbol{\theta}|\boldsymbol{Y})$ evaluated at $\boldsymbol{\theta}^{(\nu)}$;

4. Set $\boldsymbol{\theta}^{(\xi)} = \boldsymbol{\theta}^{(\nu+1)}$ and evaluate (2.10):

- If (2.10) is positive definite, then set $\boldsymbol{\theta}^{(\nu+1)}$ as the updated parameter;

- Otherwise, keep repeating steps 3 and 4 with smaller and smaller $\kappa^{(\nu)}$, until (2.10) is positive definite.

In cases where the complete-data logarithm likelihood does not allow for closed form solution of parameter estimators, RAI and MATTHEWS (1993) have found that the GEM-NR can reduce significantly the computation burden, when compared to the EM algorithm. In Appendix B, we derived all expressions (first and second order derivatives of the complete-data logarithm likelihood) to implement the GEM-NR algorithm for estimation of parameters in the MTMIM model.

## 2.5 Genome-wide score-based threshold

The genome-wide threshold in the multiple trait CIM model of JIANG and ZENG (1995) is either based on asymptotic approximation of the LRT to the chi-squared distribution or permutation. In this section we extend the score statistic of ZOU *et al.* (2004) to assess the genome-wide statistical significance level of any effect of QTL in the MTMIM model (2.1). Based on the individual and overall likelihoods, we derived all required expressions to compute the score statistic to test any effect parameter in the MIMIM model (see Appendix B).

In the genome-wide scan context of MTMIM model, the score is computed in the same fashion as in the MIM model detailed in Chapter 1, except for the dimensionality of the score function. In genome-wide scan of the MTMIM model, a putative pleiotropic QTL is assumed at every position $\lambda \in \boldsymbol{\zeta}$ and the significance level of its effects (main or epistatic effects) are tested against the null of no effects. For instance, assume a model with $m - 1$ QTL with main effects and $p$ epistatic effects between certain QTL. Assume we are scanning for a putative $m^{th}$ QTL. Let $l = \lambda$ denotes the testing position of the putative QTL coming into the model. Let

$\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \cdots, \lambda_{(m-1)}, l)$ be the current positions of all $m$ QTL in the model. Let $\boldsymbol{\theta}_m = \boldsymbol{\beta}'_m$ be a $T$ by 1 vector of effects for the new QTL coming into the model, and let $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_{m-1}, \boldsymbol{\theta}_m, \boldsymbol{\theta}_{m+1}, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}', vect(\boldsymbol{\Sigma}_e))'$ be a column vector of all parameters in the model, where $\boldsymbol{\theta}_b = \boldsymbol{\beta}'_b$ for $1 \le b \le m$ and $\boldsymbol{\theta}_b = \boldsymbol{w}'_b$ for $m < b \le s = m + p$. Let $\boldsymbol{\eta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_{m-1}, \boldsymbol{\theta}_{m+1}, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}', vect(\boldsymbol{\Sigma}_e))'$ be the column vector of nuisance parameters. Then the hypotheses $H_0 : \boldsymbol{\theta}_m = \boldsymbol{0}$ and $H_1 : \boldsymbol{\theta}_m \ne \boldsymbol{0}$ are tested at every position $l$ in the genome by means of the LRT. The genomic position with the maximum LRT over all $l$ is assessed for the presence of a QTL with the score-based threshold.

The score-based threshold is computed in the same fashion as for the MIM model. In order to maintain the variances of the re-sampled score and score statistic equal, we multiply $\hat{\boldsymbol{U}}_i$ by random variables $z_i$ from the normal distribution with mean zero and unit variance (N(0,1)). The the steps of the re-sampling score algorithm in the MTMIM model are:

1. generate $n$ independent normal variables $z_i$ $(i = 1, 2, \cdots, n)$ from N(0,1);

2. for each $l$, compute $\hat{\boldsymbol{U}}^*(l) = \sum_{i=1}^{n} \hat{\boldsymbol{U}}_i(l)z_i$, $S^*(l) = \hat{\boldsymbol{U}}^{*\prime}(l)\hat{\boldsymbol{V}}^{-1}(l)\hat{\boldsymbol{U}}^*(l)$. Then, compute $S^* = \max_{l \in \boldsymbol{\zeta}}\{S^*(l)\}$;

3. repeat *steps* 1 and 2 many times, say $N$ times (re-sampling), to obtain a sequence $(S_1^*, S_2^*, \cdots, S_N^*)$;

4. the score-based threshold for a given significance $\alpha$-level is the $100(1-\alpha)$ percentile of the ascending ordered values $(S_{(1)}^*, S_{(2)}^*, \cdots, S_{(N)}^*)$.

where, $\hat{\boldsymbol{V}}(l) = \sum_{i=1}^{n} \hat{\boldsymbol{U}}_i(l)\hat{\boldsymbol{U}}'_i(l)$, and $\hat{\boldsymbol{U}}_i(l)$ is:

$$
\hat{\boldsymbol{U}}_i(l) = \frac{\partial \ell_i\,(\boldsymbol{\theta}_m, \boldsymbol{\eta})}{\partial \boldsymbol{\theta}_m}\Bigg|_{(\boldsymbol{\theta}_m = 0, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})}
$$
$$
- \frac{\partial \ell\,(\boldsymbol{\theta}_m, \boldsymbol{\eta})}{\partial \boldsymbol{\theta}_m \partial \boldsymbol{\eta}'}\Bigg|_{(\boldsymbol{\theta}_m = 0, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})} \left( \frac{\partial \ell\,(\boldsymbol{\theta}_m, \boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}'}\Bigg|_{(\boldsymbol{\theta}_m = 0, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})} \right)^{-1} \frac{\partial \ell_i\,(\boldsymbol{\theta}_m, \boldsymbol{\eta})}{\partial \boldsymbol{\eta}}\Bigg|_{(\boldsymbol{\theta}_m = 0, \boldsymbol{\eta} = \tilde{\boldsymbol{\eta}})}
$$

where, $\tilde{\boldsymbol{\eta}}$ is the MLE of $\boldsymbol{\eta}$ under $H_0$ (see Appendix B for a detailed derivation of first and second order derivatives of the likelihood function).

## 2.6   Model selection and model optimization

The search for genomic positions associated with changes in the phenotype between groups of subjects due to QTL consists on identifying a subset of regions in the genome for which the effects of QTL are significantly different from zero. BROMAN and SPEED (2002) elaborated the problem of finding such a subset of regions in the genome, as the one of model selection, for which there exits many tools available in the vast literature of subset selection. However, in QTL studies the identification of a reasonable model, which maximizes the correct number of QTL while controlling the rate of false discovery, is predominant over the identification of models with the smallest prediction errors, which is the major criterion for model selection in the literature (BROMAN and SPEED, 2002).

The score-based threshold can be used as a criterion to build and refine models with many QTL. Starting with a model with no QTL effect we can select significant putative QTL based on the score-based threshold, add them to the model, and then further refine the model, by including to or excluding from the MTMIM model effects of QTL. We propose an algorithm, analogue to the algorithm described in ZENG *et al.* (1999), to build an initial model and to refine it upon using the score-based threshold criterion. *Forward selection*– assuming that model (2.1) starts with no QTL, one QTL is added at each step of the forward selection. In the $m^{th}$ step of the forward selection, we assume a putative pleiotropic QTL at every position $l \in \boldsymbol{\zeta}$, but avoiding positions within 5 cM neighboring regions of the $m-1$ QTL already identified in previous steps of the forward selection and compute the MLE of all parameters in the MTMIM model with QTL at positions $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \cdots, \lambda_{m-1}, l)$. For each position $l$, we compute the LRT statistic to test the null hypothesis $H_0 : (\beta_{1m}, \beta_{2m}, \cdots, \beta_{Tm})' = (0, 0, \cdots, 0)'$ for putative $m^{th}$ QTL versus $H_1 : (\beta_{1m}, \beta_{2m}, \cdots, \beta_{Tm})' \neq (0, 0, \cdots, 0)'$. A putative QTL

at the maximum LRT statistic over all positions $l$ is added to the model if the LRT statistic is larger than the score-based threshold. Next, the effect of the selected QTL on each trait is tested individually against the null of no effect using the LRT and critical value from a chis-squared probability distribution function with one degree of freedom and pre-specified corrected experiment-wise error rate $\alpha$, i.e., when $T$ traits are analyzed jointly, the corrected significance level (Bonferroni correction) to test each effect of the $m^{th}$ QTL is $\alpha_c = \alpha/T$. Finally, any non-significant effect of the $m^{th}$ QTL is removed from the model, ending the $m^{th}$ step of the forward selection. The forward selection continues until no maximum LRT statistic exceeds the score-based threshold. *Model optimization*: in turns, we update the positions of all QTL selected in the forward selection. We pick a QTL in the model, and hold the other QTL in the fixed at the positions that they were found in the forward selection. The effects of the picked QTL are then removed from the model and a new search for a QTL is done within the region delimited by its two neighboring QTL, avoiding 5 cM from the neighbor QTL (the search is performed until the end of the chromosome if no neighbor QTL is found on either side of the picked QTL). The new position of the picked QTL is set to the position of the maximum LRT statistic within the searched region and all parameters in the model are updated. This procedure is repeated until all positions of QTL in the MTMIM model are updated.

## 2.7   Testing for pleiotropy versus close linkage

Although testing for pleiotropy versus close linkage is part of model selection, we preferred to separate it from the model selection because, generally, this test is performed at the end of the model selection procedure, when the final model is almost fitted.

When models are nested, the critical value to assess the strength of the LRT is straightforward in the sense that the asymptotic distribution of the LRT is known to be $\chi^2$ with degrees of freedom equal to the difference between the number of

parameters in the full and reduced model. However, the pleiotropic and close linkage models may not be nested (for instance, models (1.8) and (1.9)), which then requires some correction of for LRT (VUONG, 1989; KAPETANIOS and WEEKS, 2003). The parametric bootstrap method of KAPETANIOS and WEEKS (2003) is an alternative for computing the empirical distribution of the LRT statistic when models are not nested. In the multiple trait model of KNOTT and HALEY (2000) a parametric bootstrap is advocated to estimate the threshold for the LRT. In recognizing the testing of pleiotropy versus close linkage as one of model selection, alternative criteria, such as BIC, AIC, AICc, bootstrap, and CV may be used to choose among competing models. In the next chapter we evaluate the performance of the AICc and LRT, using simulation.

## 2.8 Concluding remarks

A novel statistical method for multiple trait multiple interval mapping (MTMIM) of QTL from inbred line crosses was proposed. Our model belongs to the class of *seemingly unrelated regression* models (ZELLNER, 1962). We also proposed a novel method for estimation of genome-wide threshold to assess the significance level of effect of putative QTL in the MTMIM model. The method of genome-wide threshold estimation is based on the score-based framework of ZOU *et al.* (2004). Our MTMIM model has the advantage of allowing for mapping QTL with effects only on a subset of traits under analysis, while taking advantage of correlation between traits.

Our method provides a comprehensive framework for QTL inference on multiple trait and the score-based threshold serves as an essential and elegant tool for computing significance level of effects of putative QTL in genome-wide scan. Therefore, our procedure overcomes the drawbacks of permutation test as proposed in the multiple trait least square method of KNOTT and HALEY (2000), and lack of criterion for assessing significance level of effects of putative QTL in the multiple trait Bayesian method of BANERJEE *et al.* (2008).

# 3

# Evaluation of the MTMIM model by simulation study and experimental data analysis

In Chapter 1, we reviewed the current status of research on statistical methods for mapping multiple QTL in single and multiple complex traits within the maximum likelihood and Bayesian frameworks. In Chapter 2, we proposed a MTMIM model for QTL inference from inbred line crosses. We also proposed a novel method for estimation of genome-wide threshold to assess the significance level of effect of putative QTL in the MTMIM model. The method of genome-wide threshold estimation is based on the score statistic framework of ZOU *et al.* (2004). Our MTMIM model has the advantage of allowing for mapping QTL with effects only on a subset of traits under analysis, while taking advantage of correlation between traits. We ended Chapter 2 with a brief description of some alternative criteria for testing pleiotropy versus close linkage models. In this chapter, we implement our MTMIM model and score-based threshold method and evaluate them with several simulated data sets. More specifically, we evaluate type I error (Section 3.1), model fitting (Section 3.2), and pleiotropic versus close linkage model testing criteria (Section 3.3). We end this

chapter with an analysis of data from an experiment with *Drosophila* (Section 3.4).

## 3.1 Genome-wide type I error

In this section, we use simulation to evaluate the proportion of falsely discovered QTL (type I error) in the analysis of data without QTL effects. The LRT statistic is used for hypothesis testing and the score-based threshold is used as the criterion to assess significance level of effect of QTL in genome-wide scan. Each replicate has six chromosomes, each with 9 markers evenly spaced 10 cM apart from each other, 300 subjects, and three traits with parameters shown in Table 3.1. In the genome-wide scan a putative pleiotropic QTL with main effects on all traits, $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3)'$, was assumed at each 1 cM in the genome for the alternative hypothesis. The effects of QTL were tested against the simulated null hypothesis of no effects, $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3)' = (0, 0, 0)'$. For each position in the genome, we re-sampled the score statistic 1000 times to obtain the genome-wide score-based threshold. One thousand replicates were analyzed with the MTMIM model and the results are shown in Figure 3.1A. The results show clearly an excellent agreement between estimated type I error and nominal level in the range of 1 to 15%.

The empirical distributions of genome-wide score- and permutation-based thresholds in interval mapping analysis were very similar to each other (Figure 3.1B). In a fixed position in the genome, the score- and permutation-based distributions approximate that of a $\chi_3^2$ (Figures 3.1C and 3.1D) because the hypothesis being tested involves three parameters simultaneously, namely, $H_0 : (\beta_1, \beta_2, \beta_3)' = (0, 0, 0)'$ versus $H_1 : (\beta_1, \beta_2, \beta_3)' \neq (0, 0, 0)'$.

Table 3.1: Genetic architecture of traits (T1, T2 and T3) dictated by QTL (Q1, Q2, $\cdots$, Q5).

| Scenario[a] | | $h^{2}$[&] | $\mu$[b] | \multicolumn{5}{c}{Effects of each QTL ($\beta$)[c]} | | | | \multicolumn{3}{c}{$\Sigma_e$[$]} | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Q1 | Q2 | Q3 | Q4 | Q5 | T1 | T2 | T3 |
| S0 | T1 | 0 | 30 | 0 | 0 | 0 | 0 | 0 | 1 | 0.2 | 0 |
| | T2 | 0 | 35 | 0 | 0 | 0 | 0 | 0 | 0.2 | 1 | -0.2 |
| | T3 | 0 | 30 | 0 | 0 | 0 | 0 | 0 | 0 | -0.2 | 1 |
| SI | T1 | 25 | 30 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 1 | 0.2 | 0 |
| | T2 | 25 | 35 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.2 | 1 | -0.2 |
| | T3 | 25 | 30 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0 | -0.2 | 1 |
| | Chr. | – | – | 1 | 2 | 3 | 5 | 6 | – | – | – |
| | Position[d] | – | – | 23 | 15 | 45 | 67 | 53 | – | – | – |
| SII | T1 | 25 | 30 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 1 | 0.2 | 0 |
| | T2 | 18 | 35 | 0 | 0.54 | 0.54 | 0.54 | 0 | 0.2 | 1 | -0.2 |
| | T3 | 5 | 30 | 0 | 0 | 0.46 | 0 | 0 | 0 | -0.2 | 1 |
| | Chr. | – | – | 1 | 2 | 3 | 5 | 6 | – | – | – |
| | Position | – | – | 23 | 15 | 45 | 67 | 53 | – | – | – |
| SIII | T1 | 18 | 30 | 0.54 | 0 | 0.54 | 0 | 0.54 | 1 | 0.2 | – |
| | T2 | 18 | 35 | 0 | 0.54 | 0.54 | 0.54 | 0 | 0.2 | 1 | – |
| | Chr. | – | – | 1 | 1 | 3 | 6 | 6 | – | – | – |
| | Position | – | – | 23 | 33 | 45 | 38 | 53 | – | – | – |

[a] Scenario S0 is for type I error evaluation. Scenarios SI, SII and SIII are for model fitting evaluations.

[b] General mean of each trait.

[c] Main effect of QTL. The percentage of phenotypic variation of each trait due to each QTL is 5%.

[d] Position, in cM, of the QTL from the leftmost marker in the chromosome (Chr.).

[&] Heritability (%) due to all QTL affecting the trait.

[$] Residual variance-covariance matrix.

Figure 3.1: (A) Estimated and expected type I error, in percentage, of LRT when using the genome-wide score-based threshold to assess significance level of putative QTL in genome-wide scan of 1000 replicates. (B) Empirical distributions of permutation- and score-based genome-wide thresholds for one replicate simulated without QTL effects. (C) and (D) show the quantile-quantile plots of a $\chi_3^2$ distribution versus the permutation- and score-based thresholds values at a fixed position in the genome, respectively. Both (C) and (D) are results from one replicate simulated without QTL effects.

## 3.2 Model fit evaluations

In this section, we use simulation to evaluate the overall performance of our MT-MIM model and score-based threshold as the criterion to assess the significance level of effects of QTL in the genome-wide scan. We examined the performance of the MTMIM in three different scenarios (SI, SII and SIII), each evaluated with $B = 500$ replicates and sample size of 300 subjects. Each replicate was simulated with 6 chromosomes, each with 9 markers evenly spaced 10 cM apart from each other. The genetic architecture of each scenario is described with details in Table 3.1. For each replicate we build a MTMIM model using our proposed forward selection and model optimization procedure. The genome was partitioned at 1-cM grid for genome-wide scan. For the sake of comparison, we also build a MIM model for each trait in each replicate with our proposed forward selection and model optimization procedure. For every position in the genome, the score statistic was re-sampled 800 times for the purpose of genome-wide score-based threshold estimation.

Let's re-iterate our forward selection and model optimization procedure. *Forward selection–* assuming that model (2.1) starts with no QTL, one QTL is added at each step of the forward selection. In the $m^{th}$ step of the forward selection, we assume a putative pleiotropic QTL at every position $l \in \boldsymbol{\zeta}$, but avoiding positions within 5 cM neighboring regions of the $m - 1$ QTL already mapped in previous steps of the forward selection and compute the MLE of all parameters in the MTMIM model with QTL at positions $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \cdots, \lambda_{m-1}, l)$. For each position $l$, we compute the LRT statistic to test the null hypothesis $H_0 : (\beta_{1m}, \beta_{2m}, \cdots, \beta_{Tm})' = (0, 0, \cdots, 0)'$ for putative $m^{th}$ QTL versus $H_1 : (\beta_{1m}, \beta_{2m}, \cdots, \beta_{Tm})' \neq (0, 0, \cdots, 0)'$. A putative QTL at the maximum LRT statistic over all positions $l$ is added to the model if the LRT statistic is larger than the score-based threshold. Next, the effect of the selected QTL on each trait is tested individually against the null of no effect using the LRT and critical value from a chis-squared probability distribution function with one degree of freedom and pre-specified corrected experiment-wise error rate $\alpha$, i.e., when $T$ traits

are analyzed jointly, the corrected significance level (Bonferroni correction) to test each effect of the $m^{th}$ QTL is $\alpha_c = \alpha/T$. Finally, any non-significant effect of the $m^{th}$ QTL is removed from the model, ending the $m^{th}$ step of the forward selection. The forward selection continues until no maximum LRT statistic exceeds the score-based threshold. *Model optimization*: in turns, we update the positions of all QTL selected in the forward selection. We pick a QTL in the model, and hold the other QTL fixed at the positions that they were found in the forward selection. The effects of the picked QTL are then removed from the model and a new search for a QTL is done within the region delimited by its two neighboring QTL, avoiding 5 cM from the neighbor QTL (the search is performed until the end of the chromosome if no neighbor QTL is found on either side of the picked QTL). The new position of the picked QTL is set to the position of the maximum LRT statistic within the searched region and all parameters in the model are updated. This procedure is repeated until all positions of QTL in the MTMIM model are updated.

The general goal of each simulated scenario is: SI– with a basic and favorable situation, we want to evaluate basic properties of the MTMIM model, and check whether it would deserve further investigations; SII– with a mixture of QTL affecting one, two and three traits, we want to evaluate how well the MTMIM handles QTL with effects on only a subset of traits under analysis; SIII– with presence of close linked non-pleiotropic QTL and a pleiotropic QTL, we want to evaluate the MTMIM model under more complex genetic architecture. In SIII, we build a MTMIM model for each replicate using the forward selection without testing for pleiotropic versus close linkage models. Each model built in the forward selection was then refined with a follow-up test of pleiotropy versus close linkage. The pleiotropy versus close linkage test was carried out for every pleiotropic QTL in the MTMIM model built in the forward selection.

Scenario SI depicts a trivial situation where all QTL are independent and pleiotropic to all traits. The profile of the LRT and score statistics for one replicate is shown in Figure 3.2. The result shows the similarity between the two statistics along the

genome. In the neighboring positions of QTL the score statistic tends to produce higher values.



Figure 3.2: Profile of score and LRT statistics in MTMIM model for one replicate with three traits. The vertical green lines separate the six chromosomes, and the ticks in the horizontal axis represent the positions of genetic markers.

We evaluated our MTMIM model under three genome-wide significance levels, 1, 5 and 10%. For each replicate, all QTL selected in the forward selection are defined as mapped QTL. We summarize the performance of our method with measures that are function of the LOD-$d$ support interval ($d = 1, 1.5,$ and 2) of mapped QTL. The LOD-$d$ support interval of a mapped QTL is a continuous genomic region that includes the position of the mapped QTL and all positions on its left and right sides with LOD values larger than or equal to the LOD value at the position of the mapped QTL after subtraction of a positive constant $d$ (LANDER and BOTSTEIN, 1989). Let $Q_r$, for $r \in \{1, 2, \cdots, m = 5\}$, be a simulated QTL. A simulated QTL is defined as being **paired** with a mapped QTL if the simulated and mapped QTL are near nearby. A mapped QTL is defined as being **matched** to a paired QTL if the LOD-$d$ support interval of the mapped QTL includes the paired QTL. A mapped QTL is defined as **mismatched** if it is not matched. A simulated QTL $Q_r$ is defined as **identified** if it has a matched QTL. For each simulated $Q_r$ and for each $d$, let $\Omega_{Q_r,d}$ be the set of replicates for which $Q_r$ is identified. We define $|\Omega_{Q_r,d}|$ as the number of elements in $\Omega_{Q_r,d}$. In what follows, we define our measures of model fit.

False discovery rate per replicate (FDR$_b$):

$$FDR_b(d) = \frac{\text{number of mismatched QTL in replicate } b}{\text{total of mapped QTL in replicate } b}$$

*FDR* over all replicates:

$$FDR(d) = \frac{1}{B} \sum_{b=1}^{B} FDR_b(d)$$

*Power* to identify $Q_r$:

$$Power(Q_r, d) = \frac{|\Omega_{Q_r,d}|}{B}$$

*Coverage of LOD-d support interval* of $Q_r$:

$$C(Q_r, d) = \frac{|\Omega_{Q_r,d}|}{\text{number of replicates for which } Q_r \text{ is paired with a mapped QTL}}$$

*Mean of length of LOD-d support interval* of $Q_r$: average of lengths of LOD-$d$ support interval of $Q_r$ over replicates in $\Omega_{Q_r,d}$.

*Mean of effect* of $Q_r$: average of effects of $Q_r$ over replicates in $\Omega_{Q_r,d}$.

*Mean of position* of $Q_r$:average of positions of $Q_r$ over replicates in $\Omega_{Q_r,d}$.

We also summarize results in terms of: *model size*– number of mapped QTL; *mean of score-based threshold*– average of score-based threshold to add a QTL into a model as the number of mapped QTL grew larger in the forward selection. For each model size, only replicates with that specific model size were used to compute the mean of score-based threshold.

These summary statistics have been proposed by C. Laurie, S. Wang, L. A. Carlini-Garcia and Z-B. Zeng (unpublished). Additionally, we evaluate the accuracy of the MTMIM model in estimating the genotypic variance-covariance matrix. The estimators of the genotypic variance ($\sigma_{g_t}^2$) of trait $t$, the genotypic covariance ($\sigma_{g_{tt'}}$) and correlation ($\rho_{g_{tt'}}$) between traits $t_1$ and $t_2$ (ZENG *et al.*, 1999; MAIA, 2007, pages 109-110) are:

$$\hat{\sigma}_{g_t}^2 = \sum_{r=1}^{m+p} \sum_{u=1}^{m+p} \{\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \hat{\pi}_{ij}(\boldsymbol{Z}_{[r,j]} - \bar{Z}_r)(\boldsymbol{Z}_{[u,j]} - \bar{Z}_u)\hat{\boldsymbol{\mathcal{B}}}_{[t,r]}\hat{\boldsymbol{\mathcal{B}}}_{[t,u]}\}$$

$$\hat{\sigma}_{g_{t_1 t_2}} = \sum_{r=1}^{m+p} \sum_{u=1}^{m+p} \{\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \hat{\pi}_{ij}(\boldsymbol{Z}_{[r,j]} - \bar{Z}_r)(\boldsymbol{Z}_{[u,j]} - \bar{Z}_u)\hat{\boldsymbol{\mathcal{B}}}_{[t_1,r]}\hat{\boldsymbol{\mathcal{B}}}_{[t_2,u]}\}$$

$$\hat{\rho}_{g_{t_1 t_2}} = \frac{\hat{\sigma}_{g_{t_1 t_2}}}{\sqrt{\hat{\sigma}_{g_{t_1}}^2 \hat{\sigma}_{g_{t_2}}^2}}$$

where, $\bar{Z}_r = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \hat{\pi}_{ij}\boldsymbol{Z}_{[r,j]}$, and $\hat{\pi}_{ij}$ and $\hat{\boldsymbol{\mathcal{B}}}$ are the MLE of $\pi_{ij}$ and $\boldsymbol{\mathcal{B}}$, respectively.

In what follows, we describe the results of each summary statistic separately.

*Mean of score-based threshold* (Table 3.2): a trend clearly noticed is the almost constance of score-threshold values across models with different number of QTL, for any given genome-wide significance level. Celilia et al. (unpublished) have found the same behavior in the MIM model.

Table 3.2: Mean of score-based threshold to add a QTL into a MTMIM model, when $m$ QTL are already in the model.

| Scenario | Level[1] | $m$ | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| | 1% | 19.5 | 19.6 | 19.7 | 19.9 | 19.8 | 19.8 | – |
| SI | 5% | 16.0 | 16.1 | 16.2 | 16.3 | 16.3 | 16.3 | – |
| | 10% | 14.4 | 14.4 | 14.5 | 14.6 | 14.7 | 14.7 | – |
| | 1% | 19.6 | 19.7 | 19.6 | 19.7 | 19.7 | 19.7 | 20.9 |
| SII | 5% | 16.1 | 16.1 | 16.2 | 16.2 | 16.2 | 16.2 | 16.1 |
| | 10% | 14.4 | 14.5 | 14.5 | 14.5 | 14.6 | 14.6 | 14.5 |
| | 1% | 16.9 | 17.1 | 17.1 | 17.2 | 17.1 | – | – |
| SIII | 5% | 13.6 | 13.7 | 13.7 | 13.8 | 13.7 | 13.8 | – |
| | 10% | 12.0 | 12.1 | 12.2 | 12.2 | 12.2 | 12.1 | – |

[1] Genome-wide significance level.

*Model size* (Tables 3.3, 3.4 and 3.5): as expected, the number of QTL in the MTMIM model of scenario SI (Table 3.3) is closer to the simulated parameter (5 QTL) when compared to scenario SII (Table 3.4), for any genome-wide significance level. While a QTL in both scenarios has to exceed very similar thresholds to be declared significant in the forward selection (Table 3.2), the number of traits affected by a QTL is rather different in the two scenarios. In scenario SI all QTL have effect on all traits, while in scenario SII a QTL may have effect either on one, two or three

traits. Therefore, model over-parametrization is making the detection of QTL with effects on one and two traits in scenario SII more difficult. Similar argument carries over to scenario SIII (Table 3.5), where most of the QTL have effects on one trait only.

Our results show that in general the number of QTL mapped in MTMIM model is closer to the simulated (5 QTL) than in MIM model.

Table 3.3: Number of mapped QTL (model size) in MIM and MTMIM models in scenario SI.

| Analysis (trait) | Level[1] | Number of QTL | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| MIM (T1) | 1% | 6 | 27 | 69 | 133 | 167 | 98 | 0 | 0 |
| | 5% | 0 | 0 | 17 | 56 | 196 | 223 | 8 | 0 |
| | 10% | 0 | 0 | 7 | 36 | 136 | 296 | 24 | 1 |
| MIM (T2) | 1% | 4 | 17 | 77 | 153 | 161 | 88 | 0 | 0 |
| | 5% | 1 | 2 | 74 | 193 | 213 | 5 | 0 | 0 |
| | 10% | 0 | 0 | 4 | 44 | 150 | 284 | 17 | 1 |
| MIM (T3) | 1% | 2 | 29 | 80 | 147 | 149 | 93 | 0 | 0 |
| | 5% | 0 | 5 | 21 | 74 | 161 | 231 | 8 | 0 |
| | 10% | 0 | 1 | 8 | 37 | 142 | 291 | 20 | 1 |
| MTMIM (T1,T2,T3) | 1% | 0 | 0 | 0 | 0 | 7 | 487 | 4 | 2 |
| | 5% | 0 | 0 | 0 | 0 | 1 | 472 | 24 | 3 |
| | 10% | 0 | 0 | 0 | 0 | 0 | 439 | 56 | 5 |

[1] Genome-wide significance level.

Table 3.4: Number of mapped QTL (model size) in MIM and MTMIM models in scenario SII.

| Analysis (trait) | Level[1] | Number of QTL | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| MIM (T1) | 1% | 13 | 29 | 73 | 145 | 163 | 75 | 1 | 1 |
| | 5% | 6 | 2 | 15 | 86 | 166 | 215 | 9 | 1 |
| | 10% | 7 | 1 | 7 | 48 | 140 | 274 | 20 | 3 |
| MIM (T2) | 1% | 5 | 64 | 180 | 247 | 4 | 0 | 0 | 0 |
| | 5% | 1 | 14 | 128 | 344 | 13 | 0 | 0 | 0 |
| | 10% | 0 | 4 | 86 | 375 | 34 | 1 | 0 | 0 |
| MIM (T3) | 1% | 215 | 281 | 3 | 0 | 0 | 0 | 0 | 0 |
| | 5% | 120 | 363 | 16 | 1 | 0 | 0 | 0 | 0 |
| | 10% | 83 | 378 | 34 | 5 | 0 | 0 | 0 | 0 |
| MTMIM (T1,T2,T3) | 1% | 0 | 2 | 30 | 108 | 204 | 152 | 4 | 0 |
| | 5% | 0 | 0 | 4 | 49 | 179 | 253 | 13 | 2 |
| | 10% | 0 | 0 | 0 | 33 | 130 | 310 | 23 | 4 |

[1] Genome-wide significance level.

Table 3.5: Number of mapped QTL (model size) in MTM and MTMIM models in scenario SIII.

| Analysis (trait) | Level[1] | Number of QTL | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| MIM (T1) | 1% | 7 | 59 | 178 | 255 | 1 | 0 | 0 | 0 |
| | 5% | 2 | 14 | 111 | 359 | 14 | 0 | 0 | 0 |
| | 10% | 0 | 7 | 81 | 386 | 26 | 0 | 0 | 0 |
| MIM (T2) | 1% | 5 | 50 | 183 | 261 | 1 | 0 | 0 | 0 |
| | 5% | 0 | 11 | 92 | 383 | 14 | 0 | 0 | 0 |
| | 10% | 0 | 4 | 64 | 408 | 23 | 1 | 0 | 0 |
| MTMIM (T1, T2) | 1% | 1 | 3 | 25 | 159 | 213 | 96 | 3 | 0 |
| | 5% | 0 | 1 | 9 | 149 | 215 | 119 | 7 | 0 |
| | 10% | 0 | 0 | 2 | 135 | 217 | 131 | 15 | 0 |

[1] Genome-wide significance level.

*FDR* (Tables 3.6, 3.7 and 3.8): FDR is a very import measure of quality control in statistical analyses. However, FDR is not feasibly estimated in analysis of data from traditional QTL experiments, due to the low discovery rate of putative QTL in such experiments. Nevertheless, in simulation experiments we are able to estimate FDR because we can replicate the experiment many times. We estimate FDR when varying the genome-wide significance levels (1, 5, and 10%) and LOD-$d$ support interval levels ($d$=1, 1.5 and 2). While FDR is expected to increase with increments in genome-wide significance level, our results show that for a fixed LOD-$d$ level FDR changed few with increments in genome-wide significance levels, in both MIM and MTMIM models. Regarding changes in LOD-$d$ level, our results show that FDR and LOD-$d$ are negatively correlated, as expected. Higher levels of LOD-$d$ ultimately translate into wider LOD-$d$ support intervals. Therefore, increasing chances of capturing the true position of QTL. FDR in MIM and MTMIM models were very similar, except MIM model of trait T3 of scenario SII (Table 3.8), which was simulated with only one QTL of small effect (heritability of 5%).

Table 3.6: False discovery rate (FDR) estimates in MIM and MTMIM models in scenario SI.

| Analysis (trait) | $d^1$ | Genome-wide significance level | | |
|---|---|---|---|---|
| | | 1% | 5% | 10% |
| MIM (T1) | 1.0 | 0.091 | 0.091 | 0.099 |
| | 1.5 | 0.039 | 0.044 | 0.053 |
| | 2.0 | 0.020 | 0.027 | 0.036 |
| MIM (T2) | 1.0 | 0.080 | 0.087 | 0.089 |
| | 1.5 | 0.039 | 0.042 | 0.047 |
| | 2.0 | 0.020 | 0.023 | 0.030 |
| MIM (T3) | 1.0 | 0.107 | 0.096 | 0.099 |
| | 1.5 | 0.038 | 0.042 | 0.049 |
| | 2.0 | 0.018 | 0.023 | 0.031 |
| MTMIM (T1, T2, T3) | 1.0 | 0.046 | 0.054 | 0.069 |
| | 1.5 | 0.019 | 0.027 | 0.040 |
| | 2.0 | 0.011 | 0.019 | 0.033 |

[1] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.7: False discovery rate (FDR) estimates in MIM and MTMIM models in scenario SII.

| Analysis (trait) | $d^1$ | Genome-wide significance level | | |
|---|---|---|---|---|
| | | 1% | 5% | 10% |
| MIM (T1) | 1.0 | 0.089 | 0.092 | 0.100 |
| | 1.5 | 0.037 | 0.043 | 0.053 |
| | 2.0 | 0.018 | 0.022 | 0.030 |
| MIM (T2) | 1.0 | 0.079 | 0.086 | 0.096 |
| | 1.5 | 0.032 | 0.041 | 0.054 |
| | 2.0 | 0.012 | 0.022 | 0.036 |
| MIM (T3) | 1.0 | 0.124 | 0.138 | 0.180 |
| | 1.5 | 0.075 | 0.090 | 0.114 |
| | 2.0 | 0.048 | 0.065 | 0.085 |
| MTMIM (T1, T2, T3) | 1.0 | 0.085 | 0.092 | 0.100 |
| | 1.5 | 0.033 | 0.041 | 0.049 |
| | 2.0 | 0.014 | 0.024 | 0.032 |

[1] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.8: False discovery rate (FDR) estimates in MIM and MTMIM models in scenario SIII.

| Analysis (trait) | $d^1$ | Genome-wide significance level | | |
|---|---|---|---|---|
| | | 1% | 5% | 10% |
| MIM (T1) | 1.0 | 0.072 | 0.079 | 0.087 |
| | 1.5 | 0.028 | 0.035 | 0.041 |
| | 2.0 | 0.014 | 0.019 | 0.023 |
| MIM (T2) | 1.0 | 0.062 | 0.070 | 0.078 |
| | 1.5 | 0.031 | 0.037 | 0.045 |
| | 2.0 | 0.012 | 0.021 | 0.028 |
| MTMIM (T1, T2) | 1.0 | 0.056 | 0.078 | 0.084 |
| | 1.5 | 0.029 | 0.052 | 0.057 |
| | 2.0 | 0.022 | 0.041 | 0.045 |

[1] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

*Power* (Tables 3.9, 3.10 and 3.11): results of MIM and MTMIM models for all three scenarios clearly show a remarkable increment in power as genome-wide significance levels grow less stringent, for any LOD-$d$ level. Based on these results as well as those that showed almost constance of FDR across genome-wide significance levels, we, hereafter, show results of 10% genome-wide significance level only. The results of all other significance levels are shown in complementary tables in Appendix D.

Results of power (10% genome-wide significance level) to identifying QTL in MT-MIM model show that QTL affecting more traits have higher chances of being identified in the forward selection. In scenario SI, which is the most favorable amongst all three scenarios, all QTL have effects on all traits. Therefore, all QTL were correctly identified most of the times, power $\geq 93\%$ (Table 3.9). In scenario SII, Q1 has effect on one trait only, Q2 on two traits, and Q3 on three traits. Power increases from Q1 to Q3 in MTMIM model (Table 3.10). Results also show that MTMIM model might have lower power than MIM model for QTL with effects on only a small subset of traits under analysis. For instance, MTMIM model has less power than MIM model to identify Q1, which affects only T1 (same pattern is seen for Q5). However, as the subset of traits affected by a QTL increases, power of MTMIM model overpasses power of MIM model, even when some traits are not affected by that QTL. For instance, Q2 affects T1 and T2, but not T3, nevertheless, MTMIM model identifies Q2 more frequently than MIM model (same pattern carries over to Q4). The increment in power as the number of traits affected by a QTL increases was also observed in scenario SIII (Table 3.11).

Table 3.9: Power (%) of QTL identification in MIM and MTMIM models when using forward selection and score-based threshold as the criterion for variable selection in scenario SI.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 63.2 | 66.8 | 68.4 | 78.2 | 82.0 | 83.6 | 82.2 | 86.6 | 88.4 |
| | Q2 | 60.6 | 63.6 | 64.2 | 77.8 | 81.8 | 82.6 | 84.2 | 87.6 | 89.0 |
| | Q3 | 62.8 | 67.4 | 68.8 | 77.2 | 81.6 | 83.4 | 82.0 | 87.2 | 88.8 |
| | Q4 | 63.0 | 66.4 | 68.0 | 78.4 | 81.8 | 83.6 | 83.4 | 87.0 | 88.8 |
| | Q5 | 63.6 | 66.8 | 68.0 | 79.2 | 83.6 | 85.0 | 81.8 | 86.4 | 88.0 |
| MIM (T2) | Q1 | 62.8 | 64.8 | 65.4 | 76.4 | 80.0 | 80.6 | 85.0 | 88.2 | 89.0 |
| | Q2 | 62.6 | 64.8 | 65.6 | 77.4 | 80.0 | 81.2 | 82.0 | 84.8 | 86.2 |
| | Q3 | 63.6 | 65.6 | 68.2 | 77.0 | 79.8 | 82.6 | 80.6 | 83.4 | 86.2 |
| | Q4 | 61.8 | 66.0 | 66.6 | 76.8 | 82.4 | 83.0 | 81.2 | 87.0 | 87.6 |
| | Q5 | 64.6 | 68.4 | 70.0 | 78.4 | 83.0 | 85.8 | 84.6 | 88.8 | 91.2 |
| MIM (T3) | Q1 | 59.4 | 65.6 | 67.0 | 75.4 | 81.4 | 83.4 | 80.2 | 86.0 | 88.0 |
| | Q2 | 60.4 | 63.2 | 64.2 | 77.0 | 80.0 | 80.8 | 83.6 | 86.6 | 87.6 |
| | Q3 | 59.8 | 65.6 | 67.2 | 73.6 | 80.4 | 82.4 | 77.6 | 84.0 | 86.0 |
| | Q4 | 61.2 | 65.4 | 67.0 | 77.0 | 80.8 | 82.2 | 84.0 | 87.8 | 89.2 |
| | Q5 | 61.0 | 65.4 | 66.8 | 79.6 | 83.0 | 84.8 | 85.0 | 88.6 | 90.8 |
| MTMIM (T1,T2,T3) | Q1 | 96.2 | 98.8 | 99.0 | 96.4 | 99.4 | 99.6 | 96.4 | 99.4 | 99.6 |
| | Q2 | 97.0 | 98.0 | 98.6 | 97.0 | 98.0 | 98.6 | 96.8 | 98.2 | 98.8 |
| | Q3 | 94.2 | 97.0 | 98.4 | 94.6 | 97.4 | 98.8 | 94.6 | 97.4 | 98.8 |
| | Q4 | 93.6 | 98.4 | 99.0 | 94.0 | 98.8 | 99.4 | 93.8 | 99.0 | 99.4 |
| | Q5 | 96.4 | 98.6 | 99.6 | 96.4 | 98.6 | 99.6 | 96.0 | 98.6 | 99.6 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amounts subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.10: Power (%) of QTL identification in MIM and MTMIM models when using forward selection and score-based threshold as the criterion for variable selection in scenario SII.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 62.0 | 65.8 | 67.8 | 76.2 | 80.2 | 83.0 | 80.4 | 84.2 | 87.4 |
| | Q2 | 57.0 | 59.8 | 61.0 | 74.6 | 78.2 | 79.2 | 78.8 | 81.8 | 83.8 |
| | Q3 | 58.8 | 63.2 | 64.4 | 75.6 | 81.2 | 82.8 | 79.5 | 85.8 | 88.0 |
| | Q4 | 60.6 | 63.4 | 64.6 | 75.6 | 78.4 | 80.2 | 80.2 | 83.4 | 85.6 |
| | Q5 | 62.0 | 65.6 | 66.4 | 77.6 | 82.0 | 83.2 | 82.4 | 87.2 | 87.8 |
| MIM (T2) | Q2 | 70.6 | 74.4 | 75.2 | 81.0 | 85.4 | 86.4 | 85.4 | 89.8 | 91.0 |
| | Q3 | 72.2 | 76.4 | 78.4 | 81.2 | 86.0 | 88.2 | 85.0 | 90.0 | 92.2 |
| | Q4 | 74.4 | 77.4 | 79.0 | 84.8 | 87.6 | 89.2 | 89.4 | 92.0 | 93.4 |
| MIM (T3) | Q3 | 50.6 | 53.4 | 55.0 | 67.2 | 70.6 | 72.4 | 73.0 | 77.8 | 79.6 |
| MTMIM (T1, T2, T3) | Q1 | 49.8 | 53.8 | 55.2 | 66.2 | 71.0 | 72.4 | 73.2 | 78.2 | 79.8 |
| | Q2 | 83.8 | 89.0 | 89.8 | 89.8 | 94.4 | 95.6 | 91.0 | 95.6 | 96.8 |
| | Q3 | 93.0 | 96.6 | 99.0 | 92.8 | 97.0 | 99.4 | 92.8 | 97.2 | 99.4 |
| | Q4 | 85.0 | 87.6 | 89.0 | 90.2 | 93.2 | 94.4 | 91.4 | 94.6 | 96.0 |
| | Q5 | 52.0 | 57.2 | 58.6 | 65.6 | 71.8 | 73.2 | 71.8 | 78.4 | 80.0 |

[a] Genome-wide significance level.
[b] 1, 1.5 and 2 are the amounts subtracted from the LOD value at QTL position to estimate LOD-$d$ support interval for the QTL.

Table 3.11: Power (%) of QTL identification in MIM and MTMIM models when using forward selection and score-based threshold as the criterion for variable selection in scenario SIII.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 66.6 | 67.6 | 67.8 | 76.0 | 77.2 | 77.4 | 78.2 | 79.6 | 79.8 |
| | Q3 | 72.2 | 75.2 | 76.8 | 83.4 | 87.0 | 89.0 | 86.4 | 90.2 | 92.4 |
| | Q5 | 67.2 | 70.2 | 70.2 | 75.2 | 78.4 | 78.4 | 78.4 | 81.6 | 81.6 |
| MIM (T2) | Q2 | 62.6 | 64.2 | 64.8 | 71.6 | 74.2 | 74.6 | 73.6 | 76.4 | 76.8 |
| | Q3 | 73.2 | 76.4 | 78.0 | 85.2 | 88.4 | 90.2 | 87.6 | 91.2 | 93.0 |
| | Q4 | 73.2 | 74.6 | 76.2 | 84.0 | 86.0 | 87.4 | 86.0 | 88.0 | 89.4 |
| MTMIM (T1, T2) | Q1 | 63.4 | 65.4 | 65.6 | 63.4 | 65.2 | 65.8 | 67.8 | 70.0 | 70.6 |
| | Q2 | 63.4 | 64.6 | 65.0 | 65.4 | 66.6 | 67.2 | 66.6 | 68.0 | 69.0 |
| | Q3 | 91.0 | 94.4 | 96.0 | 92.8 | 96.4 | 98.0 | 93.6 | 97.0 | 98.8 |
| | Q4 | 72.6 | 74.8 | 75.6 | 75.6 | 77.4 | 78.8 | 76.0 | 78.2 | 79.6 |
| | Q5 | 64.4 | 65.6 | 65.6 | 64.8 | 66.2 | 66.6 | 66.4 | 68.0 | 68.4 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amounts subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

In scenarios SII and SIII, we decomposed power of QTL identification into three non-overlapping subsets. In scenario SII, there is a subset of replicates for which a QTL affects T1 only, another subset for which a QTL affects T1 and T2 simultaneously, and finally a subset of replicates for which a QTL affects all traits simultaneously (Table 3.12). In scenario SIII, there is a subset of replicates for which a QTL affects T1 only, another subset for which a QTL effects T2 only, and finally a subset of replicates for which a QTL affects T1 and T2 simultaneously (Table 3.13). These decompositions of power allow us to separate power of MTMIM model into QTL-by-trait power. Therefore, enabling us to measure the frequency in which a non-pleiotropic QTL is mapped as a pleiotropic QTL. In scenario SII, where all QTL are independent, most of power to identifying a QTL is concentrated on the right trait affected by that QTL. For instance, in the LOD-1.5 level, 66.4 out of 78.2 power (0.85 ratio) to identifying Q1 is due to T1 only, which is the only trait that Q1 has effect on (Table 3.12). In scenario SIII, because linkage between QTL pairs Q1 and Q2, and Q3 and Q4, contribution of power to identifying a QTL due to the right trait affected by that QTL is lower than in scenario SII, thought the right trait still accounts for a large amount of power. For example, in the LOD-1.5 level, 36.8 out of 70 power (0.53 ratio) to identifying Q1 is due to T1 only, which is the only trait Q1 has effect on, and 46 out 68 (0.68 ratio) power to identifying Q5 is due to T1 only, which is the only trait Q5 has effect on (Table 3.13). Note that Q1 was mapped as a pleiotropic QTL (subset (1,1) in Table 3.13) more often than Q5, for instance, in the LOD-1.5 level, 30.4 out 70 (0.43 ratio) and 20.8 out of 68 power (0.31 ratio), respectively. Identification of Q1 as pleiotropic more often than Q5 is mainly because the distance between Q1 and Q2 is shorter than distance between Q4 and Q5, 10 and 15 cM, respectively. The smaller the distance between two non-pleiotropic QTL, the harder is to separate them in the MTMIM model. Moreover, separation of non-pleiotropic QTL is also affected by distance between genetic markers. Linkage maps with markers closely spaced are expected to help in separating non-pleiotropic QTL. On the other hand, separation of non-pleiotropic QTL in linkage maps with sparse

markers, such as the linkage map used in our simulations, is a much harder task. We return to the issue of separating non-pleiotropic QTL in the MTMIM model in section 3.3.

Table 3.12: Decomposition of total power ($P_{total}$) of QTL identification in scenario SII (Table 3.10) into QTL-by-trait power ($P_{trait}$) for 10% genome-wide significance level. Subsets (1, 0, 0), (1, 1, 0) and (1, 1, 1) contain replicates with QTL affecting T1 only, T1 and T2, and T1, T2 and T3, respectively. We show only three subsets out of 7 subsets of the full decomposition. These three subsets account for most of the space, besides they are the most interesting in scenario SII. The QTL-by-trait to the overall power ratio (RATIO=$P_{trait}/P_{total}$) is also presented.

| | | Subsets | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $d^a$ | | (1,0,0) | | | | | (1,1,0) | | | | | (1,1,1) | | | | |
| | | Q1 | Q2 | Q3 | Q4 | Q5 | Q1 | Q2 | Q3 | Q4 | Q5 | Q1 | Q2 | Q3 | Q4 | Q5 |
| 1 | $P_{trait}$ | 63.0 | 1.0 | 0.0 | 0.8 | 59.8 | 4.2 | 82.8 | 4.6 | 84.0 | 6.8 | 0.8 | 6.0 | 85.2 | 5.8 | 0.2 |
| | RATIO | 0.86 | 0.01 | 0.00 | 0.01 | 0.83 | 0.06 | 0.91 | 0.05 | 0.92 | 0.09 | 0.01 | 0.07 | 0.92 | 0.06 | 0.00 |
| 1.5 | $P_{trait}$ | 66.4 | 1.2 | 0.0 | 0.8 | 64.0 | 4.2 | 86.4 | 5.0 | 87.2 | 8.2 | 0.8 | 6.6 | 89.0 | 5.8 | 0.2 |
| | RATIO | 0.85 | 0.01 | 0.00 | 0.01 | 0.82 | 0.05 | 0.90 | 0.05 | 0.92 | 0.10 | 0.01 | 0.07 | 0.92 | 0.06 | 0.00 |
| 2 | $P_{trait}$ | 67.6 | 1.2 | 0.0 | 0.8 | 64.8 | 4.2 | 87.4 | 5.2 | 88.0 | 8.6 | 0.8 | 6.8 | 90.8 | 6.4 | 0.2 |
| | RATIO | 0.85 | 0.01 | 0.00 | 0.01 | 0.81 | 0.05 | 0.90 | 0.05 | 0.92 | 0.11 | 0.01 | 0.07 | 0.91 | 0.07 | 0.00 |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.13: Decomposition of total power ($P_{total}$) of QTL identification in scenario SIII (Table 3.11) into QTL-by-trait power ($P_{trait}$) for 10% genome-wide significance level. Subsets (1, 0), (0, 1) and (1, 1) contain replicates with QTL affecting T1 only, T2 only, and T1 and T2, respectively. The QTL-by-trait to the overall power ratio (RATIO=$P_{trait}/P_{total}$) is also presented.

| $d$[a] | | \multicolumn{5}{c}{(1,0)} | | | | | \multicolumn{5}{c}{(0,1)} | | | | | \multicolumn{5}{c}{(1,1)} | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Q1 | Q2 | Q3 | Q4 | Q5 | Q1 | Q2 | Q3 | Q4 | Q5 | Q1 | Q2 | Q3 | Q4 | Q5 |
| 1 | $P_{trait}$ | 36.2 | 2.8 | 3.8 | 1.2 | 45.0 | 2.4 | 35.4 | 3.6 | 48.8 | 1.6 | 29.2 | 28.4 | 86.2 | 26.0 | 19.8 |
| | RATIO | 0.53 | 0.04 | 0.04 | 0.02 | 0.68 | 0.04 | 0.53 | 0.04 | 0.64 | 0.02 | 0.43 | 0.43 | 0.92 | 0.34 | 0.30 |
| 1.5 | $P_{trait}$ | 36.8 | 2.8 | 3.4 | 1.0 | 46.0 | 2.8 | 36.2 | 4.0 | 49.6 | 1.2 | 30.4 | 29.0 | 89.6 | 27.6 | 20.8 |
| | RATIO | 0.53 | 0.04 | 0.04 | 0.01 | 0.68 | 0.04 | 0.53 | 0.04 | 0.63 | 0.02 | 0.43 | 0.43 | 0.92 | 0.35 | 0.31 |
| 2 | $P_{trait}$ | 37.0 | 3.0 | 3.2 | 1.2 | 46.6 | 2.8 | 36.8 | 4.4 | 50.6 | 0.8 | 30.8 | 29.2 | 91.2 | 27.8 | 21.0 |
| | RATIO | 0.52 | 0.04 | 0.03 | 0.02 | 0.68 | 0.04 | 0.53 | 0.04 | 0.64 | 0.01 | 0.44 | 0.42 | 0.92 | 0.35 | 0.31 |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

*Mean of position of QTL* (Tables 3.14, 3.15 and 3.16): results show that mean of positions of QTL in MIM and MTMIM models have no qualitative difference. There is, though, a trend of smaller variation (measured in terms of standard error of mean) of estimated positions in MTMIM than in MIM model. In the MTMIM model, there is a trend of smaller variation of estimated positions for those QTL with effects on a larger subset of traits under analysis.

Table 3.14: Mean of positions of QTL (cM) in MIM and MTMIM models in scenario SI for 10% genome-wide significance level. Standard errors are in parentheses.

| Analysis (Trait) | QTL | Position[a] | LOD-$d$ level | | |
| --- | --- | --- | --- | --- | --- |
| | | | 1 | 1.5 | 2 |
| | Q1 | 23 [1] | 23.4 (0.29) | 23.6 (0.36) | 23.6 (0.38) |
| | Q2 | 15 [2] | 15.7 (0.31) | 14.9 (0.33) | 15.3 (0.38) |
| MIM (T1) | Q3 | 45 [3] | 45.6 (0.32) | 45.3 (0.36) | 45.6 (0.38) |
| | Q4 | 67 [5] | 66.8 (0.27) | 66.5 (0.29) | 66.5 (0.29) |
| | Q5 | 53 [6] | 52.6 (0.28) | 52.4 (0.33) | 52.5 (0.35) |
| | Q1 | 23 [1] | 23.7 (0.32) | 23.9 (0.39) | 23.9 (0.39) |
| | Q2 | 15 [2] | 14.1 (0.30) | 14.1 (0.32) | 14.4 (0.36) |
| MIM (T2) | Q3 | 45 [3] | 44.8 (0.32) | 44.9 (0.36) | 45.0 (0.39) |
| | Q4 | 67 [5] | 66.8 (0.27) | 66.3 (0.32) | 66.2 (0.33) |
| | Q5 | 53 [6] | 52.4 (0.27) | 52.4 (0.33) | 52.4 (0.35) |
| | Q1 | 23 [1] | 23.7 (0.31) | 24.0 (0.36) | 24.3 (0.40) |
| | Q2 | 15 [2] | 14.2 (0.31) | 14.6 (0.35) | 14.9 (0.37) |
| MIM (T3) | Q3 | 45 [3] | 44.7 (0.30) | 44.8 (0.37) | 45.1 (0.39) |
| | Q4 | 67 [5] | 67.3 (0.28) | 67.1 (0.29) | 66.8 (0.34) |
| | Q5 | 53 [6] | 52.6 (0.31) | 52.4 (0.35) | 52.3 (0.37) |
| | Q1 | 23 [1] | 23.62 (0.13) | 23.76 (0.17) | 23.75 (0.17) |
| | Q2 | 15 [2] | 14.15 (0.16) | 14.20 (0.16) | 14.22 (0.16) |
| MTMIM | Q3 | 45 [3] | 45.50 (0.16) | 45.57 (0.17) | 45.68 (0.18) |
| (T1, T2 and T3) | Q4 | 67 [5] | 67.68 (0.15) | 67.66 (0.17) | 67.62 (0.17) |
| | Q5 | 53 [6] | 52.76 (0.14) | 52.73 (0.15) | 52.82 (0.16) |

[a] Simulated position (cM) of QTL from the leftmost genetic marker in the chromosome. The chromosome in which each QTL is located is shown in between the square brackets.

Table 3.15: Mean of positions of QTL (cM) in MIM and MTMIM models in scenario SII for 10% genome-wide significance level. Standard errors are in parentheses.

| Analysis (Trait) | QTL | Position[a] | LOD-$d$ level | | |
|---|---|---|---|---|---|
| | | | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 23 [1] | 23.4 (0.29) | 23.7 (0.31) | 23.8 (0.36) |
| | Q2 | 15 [2] | 14.4 (0.30) | 14.6 (0.31) | 14.9 (0.35) |
| | Q3 | 45 [3] | 45.4 (0.32) | 45.4 (0.38) | 45.2 (0.41) |
| | Q4 | 67 [5] | 67.1 (0.27) | 66.9 (0.29) | 66.7 (0.33) |
| | Q5 | 53 [6] | 52.9 (0.28) | 52.9 (0.33) | 52.8 (0.35) |
| MIM (T2) | Q2 | 15 [2] | 14.5 (0.27) | 14.7 (0.30) | 14.9 (0.32) |
| | Q3 | 45 [3] | 45.6 (0.30) | 45.2 (0.35) | 45.4 (0.37) |
| | Q4 | 67 [5] | 67.2 (0.26) | 67.0 (0.27) | 66.6 (0.33) |
| MIM (T3) | Q3 | 45 [3] | 44.7 (0.38) | 44.7 (0.45) | 44.8 (0.47) |
| MTMIM (T1, T2 and T3) | Q1 | 23 [1] | 23.4 (0.30) | 23.5 (0.32) | 23.5 (0.33) |
| | Q2 | 15 [2] | 14.4 (0.20) | 14.4 (0.22) | 14.5 (0.23) |
| | Q3 | 45 [3] | 44.9 (0.15) | 44.9 (0.18) | 44.9 (0.19) |
| | Q4 | 67 [5] | 67.6 (0.18) | 67.6 (0.19) | 67.5 (0.21) |
| | Q5 | 53 [6] | 52.9 (0.31) | 52.8 (0.37) | 52.9 (0.38) |

[a] Simulated position (cM) of QTL from the leftmost genetic marker in the chromosome. The chromosome in which each QTL is located is shown in between the square brackets.

Table 3.16: Mean of positions of QTL (cM) in MIM and MTMIM models in scenario SIII for 10% genome-wide significance level. Standard errors are in parentheses.

| Analysis (Trait) | QTL | Position[a] | LOD-$d$ level | | |
|---|---|---|---|---|---|
| | | | 1 | 1.5 | 2 |
| | Q1 | 23 [1] | 22.1 (0.18) | 22.1 (0.22) | 22.0 (0.23) |
| MIM (T1) | Q3 | 45 [3] | 44.3 (0.27) | 44.3 (0.32) | 44.1 (0.35) |
| | Q5 | 53 [6] | 52.4 (0.23) | 52.9 (0.26) | 52.9 (0.26) |
| | Q2 | 33 [1] | 34.9 (0.21) | 35.5 (0.28) | 35.6 (0.29) |
| | Q3 | 45 [3] | 43.8 (0.30) | 43.8 (0.32) | 43.8 (0.36) |
| MIM (T2) | Q4 | 38 [6] | 36.9 (0.21) | 36.5 (0.27) | 36.2 (0.29) |
| | Q1 | 23 [1] | 23.2 (0.23) | 23.1 (0.24) | 23.0 (0.25) |
| | Q2 | 33 [1] | 33.5 (0.23) | 33.6 (0.25) | 33.7 (0.27) |
| MTMIM | Q3 | 45 [3] | 44.4 (0.19) | 44.6 (0.22) | 44.5 (0.23) |
| (T1, T2) | Q4 | 38 [6] | 38.6 (0.21) | 38.5 (0.23) | 38.5 (0.27) |
| | Q5 | 53 [6] | 51.4 (0.25) | 51.7 (0.26) | 51.8 (0.29) |

[a] Simulated position (cM) of QTL from the leftmost genetic marker in the chromosome. The chromosome in which each QTL is located is shown in between the square brackets.

*Coverage and length of LOD-d support interval for position of QTL* (Tables 3.17, 3.18 and 3.19): results show that for any LOD-$d$ level the coverage of LOD-$d$ support interval for position of QTL are not remarkably different in MIM and MTMIM models. However, results show that on average estimates of length of LOD-$d$ support interval were always larger in MIM model. Differences in length are only marginal for QTL with effects on only a small subset of traits, but there are considerable differences in length for QTL with effects on larger subset of traits under analysis. For instance, in scenario SII (Table 3.18) Q1 affects one trait only and it has mean length of LOD-1.5 support intervals of 29.4 cM in MIM and 26.4 cM in MTMIM model. On the other hand, $Q2$ affects two traits and it has mean length of LOD-1.5 support interval of 27.7 ($T_1$) and 27.9 ($T_2$) in MIM and 21.0 cM in MTMIM model. An interesting result is that the LOD-1.5 support interval produced confidence intervals for position of QTL with approximately 95% coverage in both MIM and MTMIM models.

Table 3.17: Coverage (%) and length (cM) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SI for 10% genome-wide significance level. Standard errors are in parentheses.

| Analysis (Trait) | QTL | Coverage | | | Mean of length | | |
|---|---|---|---|---|---|---|---|
| | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 92.9 | 97.9 | 100 | 20.7 (0.43) | 28.0 (0.55) | 35.7 (0.64) |
| | Q2 | 92.3 | 96.3 | 97.6 | 21.3 (0.42) | 28.4 (0.61) | 35.5 (0.74) |
| | Q3 | 90.9 | 96.7 | 98.5 | 22.3 (0.48) | 31.1 (0.64) | 39.8 (0.78) |
| | Q4 | 92.5 | 96.5 | 98.5 | 19.9 (0.38) | 26.4 (0.54) | 33.4 (0.74) |
| | Q5 | 91.3 | 96.4 | 98.2 | 21.3 (0.43) | 29.1 (0.55) | 36.4 (0.61) |
| MIM (T2) | Q1 | 95.9 | 99.6 | 100 | 22.6 (0.45) | 30.7 (0.61) | 38.4 (0.69) |
| | Q2 | 93.4 | 96.6 | 98.2 | 21.6 (0.43) | 28.4 (0.58) | 35.5 (0.74) |
| | Q3 | 91.8 | 94.9 | 98.2 | 21.9 (0.45) | 30.4 (0.63) | 39.2 (0.77) |
| | Q4 | 90.8 | 97.3 | 98.0 | 20.2 (0.38) | 26.6 (0.51) | 33.2 (0.67) |
| | Q5 | 92.2 | 96.7 | 99.4 | 21.6 (0.43) | 29.6 (0.56) | 36.7 (0.66) |
| MIM (T3) | Q1 | 90.3 | 96.9 | 99.1 | 21.7 (0.43) | 29.9 (0.59) | 37.2 (0.68) |
| | Q2 | 94.1 | 97.5 | 98.7 | 21.3 (0.39) | 27.7 (0.53) | 34.9 (0.70) |
| | Q3 | 87.4 | 94.6 | 96.9 | 22.4 (0.44) | 31.2 (0.60) | 40.3 (0.76) |
| | Q4 | 92.7 | 96.9 | 98.5 | 19.6 (0.35) | 27.8 (0.47) | 33.3 (0.69) |
| | Q5 | 93.8 | 97.8 | 100 | 21.9 (0.41) | 29.5 (0.54) | 36.4 (0.63) |
| MTMIM (T1, T2 and T3) | Q1 | 96.4 | 99.4 | 99.6 | 12.4 (0.17) | 16.0 (0.26) | 19.1 (0.30) |
| | Q2 | 96.8 | 98.2 | 98.8 | 12.6 (0.17) | 16.0 (0.22) | 19.5 (0.27) |
| | Q3 | 94.6 | 97.4 | 98.8 | 12.5 (0.17) | 15.9 (0.22) | 19.4 (0.28) |
| | Q4 | 93.8 | 99.0 | 99.4 | 12.2 (0.17) | 15.3 (0.19) | 18.2 (0.23) |
| | Q5 | 96.0 | 98.6 | 99.6 | 12.3 (0.15) | 15.6 (0.20) | 18.7 (0.25) |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.18: Coverage (%) and length (cM) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SII for 10% genome-wide significance level. Standard errors are in parentheses.

| Analysis (Trait) | QTL | Coverage | | | Mean of length | | |
|---|---|---|---|---|---|---|---|
| | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 91.4 | 95.7 | 99.3 | 21.7 (0.42) | 29.4 (0.55) | 37.3 (0.66) |
| | Q2 | 92.2 | 95.8 | 98.1 | 21.1 (0.38) | 27.7 (0.55) | 34.9 (0.73) |
| | Q3 | 88.8 | 95.8 | 98.2 | 23.7 (0.49) | 33.0 (0.67) | 41.9 (0.81) |
| | Q4 | 92.2 | 95.8 | 98.4 | 20.2 (0.35) | 26.7 (0.51) | 35.4 (0.79) |
| | Q5 | 93.4 | 98.8 | 99.6 | 21.3 (0.43) | 28.7 (0.56) | 36.4 (0.68) |
| MIM (T2) | Q2 | 92.6 | 97.4 | 98.7 | 21.0 (0.88) | 27.9 (0.55) | 34.1 (0.67) |
| | Q3 | 90.6 | 95.9 | 98.3 | 22.3 (0.38) | 29.8 (0.56) | 39.1 (0.74) |
| | Q4 | 95.3 | 98.1 | 99.6 | 19.6 (0.33) | 26.1 (0.49) | 32.6 (0.67) |
| MIM (T3) | Q3 | 88.8 | 94.6 | 96.8 | 25.3 (0.55) | 35.3 (0.74) | 46.2 (0.88) |
| MTMIM (T1, T2 and T3) | Q1 | 89.5 | 95.6 | 97.6 | 20.0 (0.38) | 26.4 (0.47) | 33.1 (0.56) |
| | Q2 | 93.1 | 97.8 | 98.9 | 16.2 (0.25) | 21.0 (0.33) | 25.3 (0.39) |
| | Q3 | 92.8 | 97.2 | 99.4 | 13.1 (0.22) | 17.2 (0.28) | 20.7 (0.33) |
| | Q4 | 94.2 | 97.5 | 98.9 | 15.6 (0.23) | 20.3 (0.31) | 24.2 (0.39) |
| | Q5 | 89.5 | 97.8 | 99.8 | 19.7 (0.41) | 26.1 (0.51) | 32.6 (0.60) |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.19: Coverage (%) and length (cM) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SIII for 10% genome-wide significance level. Standard errors are in parentheses.

| Analysis (Trait) | QTL | Coverage | | | Mean of length | | |
|---|---|---|---|---|---|---|---|
| | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 98.0 | 99.8 | 100 | 20.2 (0.40) | 26.8 (0.53) | 34.9 (0.74) |
| | Q3 | 92.1 | 96.2 | 98.5 | 21.1 (0.39) | 29.1 (0.57) | 37.8 (0.74) |
| | Q5 | 95.4 | 99.3 | 99.3 | 20.4 (0.39) | 27.8 (0.57) | 35.9 (0.75) |
| MIM (T2) | Q2 | 94.9 | 98.5 | 98.9 | 20.6 (0.44) | 28.9 (0.67) | 37.5 (0.89) |
| | Q3 | 92.4 | 96.2 | 98.1 | 22.1 (0.45) | 29.6 (0.61) | 38.4 (0.77) |
| | Q4 | 95.6 | 97.8 | 99.3 | 19.5 (0.38) | 26.9 (0.55) | 35.8 (0.75) |
| MTMIM (T1, T2) | Q1 | 97.9 | 100 | 100 | 28.8 (1.16) | 38.4 (1.34) | 47.6 (1.40) |
| | Q2 | 100 | 100 | 100 | 28.1 (1.17) | 39.4 (1.38) | 48.7 (1.41) |
| | Q3 | 93.9 | 97.4 | 99.2 | 16.8 (0.31) | 22.9 (0.48) | 28.8 (0.65) |
| | Q4 | 97.4 | 100 | 100 | 26.7 (0.91) | 37.3 (1.15) | 47.6 (1.23) |
| | Q5 | 97.1 | 99.4 | 100 | 28.9 (1.05) | 39.9 (1.27) | 49.5 (1.34) |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

*Mean of effect of QTL* (Tables 3.20, 3.21 and 3.22): in scenario SI, results show that estimates of QTL effects in MTMIM model are overall close to the parameters, with neither clear trend of under- nor over-estimation (Table 3.20). On the other hand, results show that in scenario SII estimates of effects in MTMIM model for those QTL with effects on only a subset of traits seem to have a trend of over-estimation, and possibly a trend of under-estimation of effects for those QTL with effects on all traits under analysis (Table 3.21). For instance, effects of Q1 which affects T1 only were over-estimated, and effects of Q3 which affects all traits were under-estimated for traits T1 and T3. Results of scenario SII also demonstrate the robustness of MTMIM model in estimating the effects of QTL, whereby QTL without effects on certain traits have estimates near zero, while QTL with non-zero effects have estimates with few bias. However, the robustness of MTMIM to estimate effect of QTL with few bias is less evident in scenario SIII. For instance, note that while Q2 has zero effect on T1, its effect estimate is not close to zero, for any $LOD-d$ level. In order to understand why this bias is present in Q2 of scenario SIII, we need to understand how we match a mapped to a simulated QTL. Remember that in the forward selection we search and map pleiotropic QTL, then each mapped pleiotropic QTL is tested against the alternative hypothesis of close linked non-pleiotropic QTL at the neighboring region of the mapped pleiotropic QTL. If the pleiotropic hypothesis is not rejected, we assume the QTL is pleiotropic. Then, in order to apply our summary statistics each mapped pleiotropic QTL is matched to its closest (smallest distance) simulated QTL. It could happen that a mapped pleiotropic QTL in the neighboring region of simulated Q1 and Q2 be matched to Q2, even though the major effect of the mapped pleiotropic QTL comes from Q1. Note that when the previous situation happen, we mistakenly assign the effect of Q1 (which affects only T1) to Q2 (which presumably would not affect T1), therefore, producing biased effect of Q2 on T1. Had we used another criterion (not smallest distance) to match mapped to simulated QTL or had we found a more powerful statistics to separate close linked non-pleiotropic QTL, the "bias" in Q2 would be minimized or even absent. The same explanation of "bias" carries over to

Q4 (T1), Q1 (T2) and Q5 (T2) in scenario SIII. We quoted bias to emphasize that the bias observed in scenario SIII is not due to the MTMIM estimation *per se*, but rather due to our lack of ability to separate closed linked non-pleiotropic QTL or due to our criterion to match mapped to simulated QTL.

In MIM model the effects of all QTL were over-estimated. This phenomena is expected due to genome-wide selection, and it is known as "Beavis effect" (BEAVIS, 1998). A qualitative comparison of results show that overall the estimation of effects in MTMIM model are less biased than in MIM model.

Table 3.20: Mean of effect of QTL in MIM and MTMIM models in scenario SI for 10% genome-wide significance level. Standard errors are in parentheses.

| Trait | QTL | Parameter | MIM | | | MTMIM | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 |
| | Q1 | 0.52 | 0.57 (0.006) | 0.57 (0.006) | 0.57 (0.006) | 0.52 (0.006) | 0.51 (0.007) | 0.51 (0.007) |
| | Q2 | 0.52 | 0.56 (0.006) | 0.56 (0.006) | 0.56 (0.006) | 0.51 (0.006) | 0.51 (0.006) | 0.51 (0.006) |
| T1 | Q3 | 0.52 | 0.56 (0.006) | 0.56 (0.006) | 0.55 (0.006) | 0.52 (0.006) | 0.52 (0.006) | 0.52 (0.006) |
| | Q4 | 0.52 | 0.55 (0.006) | 0.55 (0.006) | 0.55 (0.006) | 0.51 (0.006) | 0.51 (0.006) | 0.51 (0.006) |
| | Q5 | 0.52 | 0.55 (0.006) | 0.56 (0.006) | 0.56 (0.005) | 0.52 (0.007) | 0.52 (0.007) | 0.51 (0.007) |
| | Q1 | 0.52 | 0.55 (0.007) | 0.55 (0.007) | 0.55 (0.007) | 0.50 (0.007) | 0.50 (0.007) | 0.50 (0.007) |
| | Q2 | 0.52 | 0.55 (0.006) | 0.56 (0.005) | 0.56 (0.005) | 0.51 (0.006) | 0.51 (0.006) | 0.51 (0.006) |
| T2 | Q3 | 0.52 | 0.56 (0.006) | 0.56 (0.005) | 0.56 (0.005) | 0.52 (0.006) | 0.52 (0.006) | 0.52 (0.006) |
| | Q4 | 0.52 | 0.55 (0.005) | 0.55 (0.005) | 0.55 (0.005) | 0.51 (0.006) | 0.50 (0.006) | 0.50 (0.006) |
| | Q5 | 0.52 | 0.56 (0.005) | 0.55 (0.006) | 0.56 (0.006) | 0.52 (0.007) | 0.52 (0.007) | 0.52 (0.007) |
| | Q1 | 0.52 | 0.56 (0.005) | 0.56 (0.005) | 0.56 (0.005) | 0.52 (0.006) | 0.52 (0.006) | 0.52 (0.006) |
| | Q2 | 0.52 | 0.55 (0.005) | 0.55 (0.005) | 0.55 (0.005) | 0.51 (0.007) | 0.51 (0.007) | 0.51 (0.007) |
| T3 | Q3 | 0.52 | 0.55 (0.005) | 0.55 (0.005) | 0.55 (0.005) | 0.51 (0.007) | 0.51 (0.006) | 0.51 (0.006) |
| | Q4 | 0.52 | 0.55 (0.005) | 0.55 (0.005) | 0.55 (0.005) | 0.51 (0.007) | 0.52 (0.007) | 0.51 (0.007) |
| | Q5 | 0.52 | 0.56 (0.006) | 0.56 (0.006) | 0.56 (0.006) | 0.53 (0.008) | 0.53 (0.008) | 0.52 (0.008) |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.21: Mean of effect of QTL in MIM and MTMIM models in scenario SII for 10% genome-wide significance level. Standard errors are in parentheses.

| Trait | QTL | Parameter | MIM | | | MTMIM | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 |
| | Q1 | 0.52 | 0.56 (0.005) | 0.56 (0.005) | 0.56 (0.005) | 0.56 (0.005) | 0.56 (0.005) | 0.56 (0.005) |
| | Q2 | 0.52 | 0.55 (0.005) | 0.56 (0.006) | 0.56 (0.005) | 0.52 (0.007) | 0.52 (0.007) | 0.52 (0.007) |
| T1 | Q3 | 0.52 | 0.54 (0.005) | 0.54 (0.005) | 0.54 (0.005) | 0.51 (0.007) | 0.51 (0.007) | 0.50 (0.007) |
| | Q4 | 0.52 | 0.55 (0.007) | 0.55 (0.006) | 0.55 (0.006) | 0.52 (0.006) | 0.52 (0.006) | 0.52 (0.006) |
| | Q5 | 0.52 | 0.55 (0.007) | 0.55 (0.006) | 0.56 (0.006) | 0.56 (0.005) | 0.56 (0.005) | 0.56 (0.005) |
| | Q1 | 0 | - | - | - | 0.00 (0.004) | 0.00 (0.004) | 0.00 (0.004) |
| | Q2 | 0.54 | 0.57 (0.006) | 0.57 (0.006) | 0.57 (0.006) | 0.55 (0.007) | 0.54 (0.007) | 0.54 (0.006) |
| T2 | Q3 | 0.54 | 0.57 (0.005) | 0.57 (0.005) | 0.57 (0.005) | 0.54 (0.007) | 0.54 (0.007) | 0.54 (0.006) |
| | Q4 | 0.54 | 0.58 (0.005) | 0.57 (0.005) | 0.57 (0.005) | 0.55 (0.006) | 0.55 (0.006) | 0.55 (0.006) |
| | Q5 | 0 | - | - | - | 0.00 (0.005) | 0.00 (0.005) | 0.00 (0.005) |
| | Q1 | 0 | - | - | - | 0.00 (0.005) | 0.00 (0.005) | 0.00 (0.005) |
| | Q2 | 0 | - | - | - | 0.00 (0.004) | 0.01 (0.004) | 0.01 (0.003) |
| T3 | Q3 | 0.46 | 0.51 (0.006) | 0.51 (0.006) | 0.51 (0.006) | 0.44 (0.008) | 0.44 (0.008) | 0.43 (0.007) |
| | Q4 | 0 | - | - | - | 0.00 (0.003) | 0.00 (0.003) | 0.00 (0.003) |
| | Q5 | 0 | - | - | - | 0.00 (0.004) | 0.00 (0.004) | 0.00 (0.004) |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table 3.22: Mean of effect of QTL in MIM and MTMIM models in scenario SIII for 10% genome-wide significance level. Standard errors are in parentheses.

| Trait | QTL | Parameter | MIM | | | MTMIM | | |
|---|---|---|---|---|---|---|---|---|
| | | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 |
| | Q1 | 0.54 | 0.57 (0.005) | 0.57 (0.006) | 0.57 (0.006) | 0.57 (0.011) | 0.56 (0.011) | 0.57 (0.012) |
| | Q2 | 0 | - | - | - | 0.21 (0.019) | 0.20 (0.019) | 0.20 (0.019) |
| T1 | Q3 | 0.54 | 0.57 (0.005) | 0.57 (0.005) | 0.57 (0.005) | 0.53 (0.008) | 0.52 (0.008) | 0.52 (0.008) |
| | Q4 | 0 | - | - | - | 0.11 (0.016) | 0.13 (0.015) | 0.12 (0.015) |
| | Q5 | 0.54 | 0.58 (0.006) | 0.58 (0.005) | 0.58 (0.005) | 0.58 (0.009) | 0.58 (0.013) | 0.59 (0.013) |
| | Q1 | 0 | - | - | - | 0.24 (0.016) | 0.23 (0.016) | 0.24 (0.016) |
| | Q2 | 0.54 | 0.58 (0.006) | 0.58 (0.006) | 0.58 (0.006) | 0.55 (0.009) | 0.55 (0.009) | 0.55 (0.009) |
| T2 | Q3 | 0.54 | 0.57 (0.005) | 0.57 (0.005) | 0.57 (0.006) | 0.53 (0.008) | 0.54 (0.008) | 0.54 (0.007) |
| | Q4 | 0.54 | 0.58 (0.005) | 0.58 (0.006) | 0.58 (0.006) | 0.59 (0.009) | 0.60 (0.008) | 0.58 (0.010) |
| | Q5 | 0 | - | - | - | 0.09 (0.016) | 0.09 (0.015) | 0.09 (0.015) |

[a] $d$ is the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

*Mean of genotypic variance-covariance matrix*: averaged genotypic covariance matrix of traits in MTMIM model are shown in Table 3.23. Because in the MTMIM model some QTL may have undetectable effects in some traits under analysis, we evaluate consequences of estimating the genotypic variance-covariance matrix with and without removing such non-significant effects from the MTMIM model. A qualitative inspection of the results shows that genotypic variance-covariance matrix estimates were very close to parameters regardless of whether non-significant effects were excluded or not from the MTMIM model.

Table 3.23: Mean of genotypic variance-covariance matrix ($\Sigma_g$) in MTMIM model in scenarios SI, SII and SIII for 10% genome-wide significance level.

| Scenario | | Parameter ($\Sigma_g$) | | | Mean of $\Sigma_g{}^{*}$ | | | Mean of $\Sigma_g{}^{\&}$ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | T1 | T2 | T3 | T1 | T2 | T3 | T1 | T2 | T3 |
| | T1 | 0.33 | 0.33 | 0.33 | 0.35 | 0.34 | 0.33 | 0.36 | 0.34 | 0.33 |
| SI | T2 | 0.33 | 0.33 | 0.33 | 0.34 | 0.35 | 0.32 | 0.34 | 0.35 | 0.33 |
| | T3 | 0.33 | 0.33 | 0.33 | 0.33 | 0.32 | 0.35 | 0.33 | 0.33 | 0.35 |
| | T1 | 0.33 | 0.21 | 0.06 | 0.34 | 0.20 | 0.06 | 0.34 | 0.21 | 0.06 |
| SII | T2 | 0.21 | 0.22 | 0.06 | 0.21 | 0.23 | 0.06 | 0.21 | 0.23 | 0.06 |
| | T3 | 0.06 | 0.06 | 0.05 | 0.06 | 0.06 | 0.06 | 0.06 | 0.06 | 0.07 |
| SIII | T1 | 0.22 | 0.18 | – | 0.24 | 0.20 | – | 0.23 | 0.20 | – |
| | T2 | 0.18 | 0.22 | – | 0.20 | 0.24 | – | 0.20 | 0.25 | – |

In MTMIM model some QTL may have non-significant effects in some traits. The mean of genotypic variance-covariance matrix indexed by '*' and '&' were computed with and without setting non-significant effects of QTL to zero, respectively.
The standard error of means ranges from 0.001 to 0.004.

## 3.3  Pleiotropy versus close linkage

In scenario SIII, after selecting a MTMIM model in the forward selection, each mapped pleiotropic QTL was tested against the alternative of close linked non-pleiotropic QTL. In bivariate model, we performed a two-dimensional search for positions of putative close linked non-pleiotropic QTL in the neighborhood of the position of each pleiotropic QTL, as suggested in JIANG and ZENG (1995). The model with non-pleiotropic QTL that showed highest likelihood within the two-dimension search region was selected and tested against the model with pleiotropic QTL. We compared two criteria of model selection, the AICc and LRT. The critical value for the LRT at 5% significance level was obtained from a $\chi^2$ probability distribution with one degree of freedom.

Estimates of type I error and power of AICc and LRT criteria for model selection are shown in Table 3.24. Because Q3 was simulated as being pleiotropic, rejection of pleiotropic hypothesis for Q3 provides a measure of type I error. On the other hand, Q1 and Q2, and Q4 and Q5 were simulated as pairs of close linked non-pleiotropic QTL. Therefore, rejection of pleiotropic hypothesis at these QTL provides a measure of power. Under our simulation setting, the LRT performed the best as criterion of model selection. The LRT was able to keep the best balance between type I error and power (Table 3.24). Estimated frequency of rejection of pleiotropy for Q3 using the LRT agrees very well with the expected 5% nominal error, and estimated frequency of rejection of pleiotropy for Q1 and Q2 are satisfactory high taking into account that Q1 and Q2 are considerably close to each other in a linkage map with markers considerably way from each (10 cM from marker to marker). On the other hand, AICc criterion showed higher power for Q1 and Q2, but with a cost of high type I error for Q3. Moreover, because Q4 and Q5 are 15 cM apart from each other, the frequency of rejection of pleiotropy for these two QTL is higher than for Q1 and Q2, which are 10 cM apart from each other.

Table 3.24: Frequency of rejection of pleiotropy model for each pleiotropic QTL in MTMIM model in scenario SIII. The AICc and LRT criteria are compared.

| QTL | AICc | | | LRT[*] | | |
|-----|------|------|------|------|------|------|
|     | 1[&] | 1.5 | 2 | 1 | 1.5 | 2 |
| Q1 | 0.45 | 0.45 | 0.45 | 0.39 | 0.38 | 0.38 |
| Q2 | 0.45 | 0.45 | 0.44 | 0.36 | 0.36 | 0.36 |
| Q3 | 0.16 | 0.15 | 0.15 | 0.05 | 0.04 | 0.04 |
| Q4 | 0.59 | 0.54 | 0.54 | 0.45 | 0.41 | 0.41 |
| Q5 | 0.65 | 0.66 | 0.66 | 0.46 | 0.48 | 0.48 |

[*] The critical value for the LRT at 5% significance level was obtained from a $\chi^2$ probability distribution with one degree of freedom.

[&] 1, 1.5 and 2 are the amounts subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

## 3.4 Analysis of data from an experiment with fruit flies *Drosophila*

In this section, we analyze data from a cross between fruit flies *Drosophila simulus* and *D. mauritiana* with MIM and our MTMIM model. The experiment, including crosses, data acquisition, and previous analyses have been described in details in LIU *et al.* (1996) and ZENG *et al.* (2000). Briefly, males from an inbred line of *D. mauritiana* (Rob A JJ) were crossed to females from an inbred line of *D. simulus* (13w JJ) to produce an $F_1$ population of males and females. $F_1$ females were then crossed to each parental line to produce two populations of males only, mauritiana backcross (BM) and simulus backcross (BS). These two crosses were repeated twice to produce two independent populations from each backcross: BS1 (sample size n=186), BS2 (n=288), BM1 (n=192) and BM2 (n=299). Males from BM1 and BS1 were genotypically scored at 45 marker loci for which the two parental lines were homozygous for different alleles. Males from BM2 and BS2 were genotypically scored at 42 marker loci out of the same 45 marker loci that BM1 and BS1 were scored. The phenotypic values of each subject are: (1) average over both sides (left and right) of the first principal component of 100 Fourier coefficients of posterior lobe (PC1); (2) area of the posterior lobe (AREA); (3) average over both sides of the first principal component of 100 Fourier coefficients of the rescaled posterior lobe, rescaled so that it has unit area (ADJPC1); and (4) length of the foreleg tibia (TIBIA). While PC1 provides a measure of both size and shape of the posterior lobe, AREA and ADJPC1, on the other hand, provide measures of size and shape, respectively. TIBIA provides a measure of overall body size. The genotypic and phenotypic data are freely available at ftp://statgen.ncsu.edu/pub/qtlcart/data/zengetal99/.

All variables related to posterior lobe (PC1, ADJPC1 and AREA) were reported to be highly correlated between themselves in both BM1 and BS1 (correlation larger than 0.82). Therefore, suggesting the presence of pleiotropic and/or close linked QTL affecting size and shape. However, all variables related to posterior lobe were

weakly correlated with TIBIA. Because possibly sharing of most developmental process components between posterior lobe shape and size, these two traits would be tightly related most due to pleiotropic effects (Liu *et al.*, 1996). Results of CIM analyses of AREA, PC1, and ADJPC1 were very similar to each other, except for the presence of a QTL affecting both AREA and PC1 but not ADJPC1 in the interval between marker loci *Ddc* and *eve*. Therefore, this QTL affects size but no shape of the posterior lobe (Liu *et al.*, 1996).

Motivated by the fact that joint analysis of PC1 and ADJPC1 in the data from the experiment with fruit flies *Drosophila* could provide additional information to distinguish between genetic effects of QTL on size and shape of posterior lobe by: (1) testing pleiotropic versus close linked non-pleiotropic QTL, and (2) estimating the genotypic covariance between traits due to linked and pleiotropic QTL, we analyzed these two traits with our MTMIM model. In the following sections, we show results of MIM and MTMIM model of the pooled samples from BM1 and BM2 (n=192+299). Hereafter, we denoted this pooled samples as BM data. We also take advantage of this data to test our GEM-NR algorithm for maximizing the likelihood function of MTMIM model with many QTL. Using data from a genetic experiment would provide more realistic differences between performances of GEM-NR and ECM algorithms than a simulated data would provide.

The LRT profiles of MIM and MTMIM models of BM data are shown in Figure 3.3. We want to mention that although MTMIM model is expected to produce larger values of LRT than MIM model at any position in the genome (Jiang and Zeng, 1995), Figure 3.3 shows that this expectation is violated at some regions in the genome. Nevertheless, this violation is easily explained because not all positions of putative QTL in the MIM and MTMIM models coincide. Therefore, MIM models are not nested within the MTMIM model shown here. Seventeen regions in the genome showed statistical evidence of putative QTL in the MTMIM model for 10% genome-wide significance level (Figure 3.3 and Table 3.25). Out of these seventeen regions, fifteen and fourteen regions also showed statistical evidence of putative QTL

in the MIM models of PC1 and ADJPC1, respectively. Overall, the inferred genomic regions harboring putative QTL in the MTMIM model are in strong agreement with previous inferred QTL in ZENG *et al.* (2000) and LIU *et al.* (1996).

MIM models of PC1 and ADJPC1 showed statistical evidence of twelve genomic regions with statistical significant QTL affecting both traits, and five regions with statistically significant QTL affecting either one of the traits (regions 3, 6, 9 , 12 and 15 shown in Figure 3.3). We want to mention that in all these five regions, expect region 6, even for the trait with no significant effect there is still some evidence of weak effect of putative QTL, as can be notice in the LRT profiles of PC1 and ADJPC1 MIM models. Region 6, which includes marker loci *Ddc* and *eve*, was previously reported not to harbor any putative QTL with significant effect on ADJPC1 (LIU *et al.*, 1996). MTMIM model mapped these five regions either exactly or very close to their respective estimated positions in the MIM model. Moreover, the estimated effects of these five regions in the MTMIM model showed small discrepancy from those estimates in the MIM model (Table 3.25). Nevertheless, empirical results from our simulations suggest that both estimates of positions and effects of QTL in MTMIM model are more accurate than in MIM model.

Positions of QTL in regions 4, 5, 7, 10, 11, 13, 14, 16 and 17 did not coincide in the MIM models of PC1 and ADJPC1. Therefore, one could hypothesize the existence of two close linked non-pleiotropic QTL at each of these regions. In order to verify this hypothesis, we tested the hypothesis of pleiotropic QTL versus close linked non-pleiotropic QTL at each of these regions. However, on the basis of the data available the null hypothesis of pleiotropic QTL could not be rejected for any region. Thus, since PC1 contains attributes of both shape and size of posterior lobe, whereas ADJPC1 contains attributes of size only, the available data provides strong evidence that the genetic mechanisms controlling shape and size of posterior lobe are highly similar. Nevertheless, availability of linkage map with many markers closely linked and larger sample size could lead to different conclusions regarding tests of pleiotropic versus close linked non-pleiotropic QTL.

Partition of the phenotypic variance-covariance matrix between PC1 and AD-JPC1 in terms of their environmental and genotypic components, as estimated in the MTMIM model, shows that most of the phenotypic covariance between these traits is due to the genotypic component, more specifically due to pleiotropic QTL rather than close linked non-pleiotropic QTL (Table 3.25).

Figure 3.3: LRT profile of separate MIM analyses of PC1 and ADJPC1 and MTMIM analysis of PC1 and ADJPC1 (Joint) of BM data for 10% genome-wide significance level. Tick marks in the horizontal axis represent positions of genetic markers on chromosomes X, 2 and 3 (from left to right). Black squares bellow the horizontal axis indicate positions of mapped QTL in separate and joint analyses.

The possibility of fitting many traits and many QTL in MTMIM model imposes severe burden in estimation of parameters. Burden both in terms of reliability of parameter estimates (accuracy) and time to estimation (speed). The GEM-NR and ECM algorithms are two alternative approaches suitable for parameter estimation in such complex models. We evaluate these two algorithms with the BM data by fitting a MTMIM model (PC1 and ADJPC1) containing 16 QTL. The results (Figure 3.4) show a tremendous gain of GEM-NR over ECM in terms of number of iterations, 19 and 52, respectively, as well as in terms of computing user time, 8.2 and 30.6 seconds, respectively. Parameter estimates delivered in the GEM-NR and ECM were very similar (not shown).

Table 3.25: Estimates of QTL main effect on PC1 ($\hat{\beta}_1$) and ADJPC1 ($\hat{\beta}_2$) in MIM and MTMIM models of BM data for 10% genome-wide significance level. The estimated genotypic variances and covariances are also shown.

| QTL | MIM PC1 $\hat{p}^{\$}$ | MIM PC1 $\hat{\beta}_1$ | MIM ADJPC1 $\hat{p}$ | MIM ADJPC1 $\hat{\beta}_2$ | MTMIM $\hat{p}$ | MTMIM $\hat{\beta}_1$ | MTMIM $\hat{\beta}_2$ |
|---|---|---|---|---|---|---|---|
| | | | Chromosome X | | | | |
| 1 | 1 | 0.00204 | 1 | 0.01651 | 1 | 0.00212 | 0.01752 |
| 2 | 20 | 0.00184 | 20 | 0.02843 | 20 | 0.00173 | 0.02747 |
| | | | Chromosome 2 | | | | |
| 3 | – | – | 1 | 0.03042 | 1 | 0.00071 | 0.02929 |
| 4 | 14 | 0.00182 | 17 | 0.02145 | 17 | 0.00181 | 0.02197 |
| 5 | 26 | 0.00171 | 30 | 0.01407 | 29 | 0.00115 | 0.01459 |
| 6 | 71 | 0.00162 | – | – | 70 | 0.00171 | $-0.00481^{ns}$ |
| 7 | 111 | 0.00091 | 116 | 0.01467 | 116 | 0.00110 | 0.01764 |
| 8 | 144 | 0.00122 | 144 | 0.00907 | 144 | 0.00114 | 0.00820 |
| | | | Chromosome 3 | | | | |
| 9 | 5 | 0.00126 | – | – | 4 | 0.00110 | 0.01066 |
| 10 | 17 | 0.00217 | 16 | 0.05030 | 17 | 0.00218 | 0.04265 |
| 11 | 48 | 0.00328 | 44 | 0.02794 | 45 | 0.00272 | 0.02526 |
| 12 | – | – | 54 | 0.02353 | 54 | $0.00069^{ns}$ | 0.02548 |
| 13 | 82 | 0.00331 | 83 | 0.03913 | 83 | 0.00337 | 0.03939 |
| 14 | 112 | 0.00092 | 116 | 0.03236 | 115 | 0.00093 | 0.02570 |
| 15 | 129 | 0.00152 | – | – | 128 | 0.00118 | $0.00943^{ns}$ |
| 16 | 147 | 0.00073 | 146 | 0.01163 | 145 | 0.00087 | 0.00923 |
| 17 | 169 | 0.00209 | 166 | 0.02681 | 167 | 0.00210 | 0.02734 |
| #QTL | 15 | | 14 | | 17 | | |
| $\hat{\boldsymbol{\Sigma}}_p$ | $27.66\text{x}10^{-6}$ | | – | | $27.66\text{x}10^{-6}$ | | $31.79\text{x}10^{-5}$ |
| | – | | $52.26\text{x}10^{-4}$ | | $31.79\text{x}10^{-5}$ | | $52.26\text{x}10^{-4}$ |
| $\hat{\boldsymbol{\Sigma}}_g$ | $23.58\text{x}10^{-6}$ | | – | | $23.64\text{x}10^{-6}$ | | $31.42\text{x}10^{-5}$ |
| | – | | $45.30\text{x}10^{-4}$ | | $31.42\text{x}10^{-5}$ | | $45.25\text{x}10^{-4}$ |
| $\hat{\boldsymbol{\Sigma}}_e$ | $4.08\text{x}10^{-6}$ | | – | | $4.02\text{x}10^{-6}$ | | $0.37\text{x}10^{-5}$ |
| | – | | $6.96\text{x}10^{-4}$ | | $0.37\text{x}10^{-5}$ | | $7.01\text{x}10^{-4}$ |

$\$$ Position, in cM, of QTL from the leftmost genetic marker on the chromosome.

ns Non-significant main effect tested with the LRT and 5% significance level. The critical value of the LRT was obtained from the $\chi^2$ distribution function with one degree of freedom.

Figure 3.4: Comparison of performances between ECM and GEM-NR algorithms in terms of number of iterations required for convergence of the likelihood function. Both algorithms were applied to a MTMIM model with 16 QTL and traits PC1 and ADJPC1 of the BM data. The algorithms were said to have converged whenever the difference between the natural logarithm of the likelihood function of two consecutive iterations was smaller than or equal to $10^{-4}$. (A) shows the values of the natural logarithm of the likelihood function at each iteration ($\log_e(L_k)$) until convergence was reached (GEM-NR algorithm began with 5 iterations of ECM algorithm. Therefore, the first 5 iterations produced identical likelihood values in both algorithms, and because of that we omitted the first 4 iterations). (B) shows the difference between the natural logarithm of the likelihood function of two consecutive iterations until convergence was reached. In (B), y-axis was re-scaled with a logarithm of base ten to improve graphical resolution.

## 3.5   Concluding remarks

Our MTMIM model and score-based threshold proposed in Chapter 2 were evaluated through simulations. Also, we analyzed data from an experiment with *Drosophila* for the purpose of illustrating our MTMIM model and evaluating the performances of the GEM-NR and ECM algorithms developed in Chapter 2.

Results from our simulations showed many interesting features of our MTMIM model and score-based threshold. First, our score-based threshold maintained the type I error at a desired nominal level when no QTL effects were present in the simulated data. Second, discovery of spurious QTL (false discovery rate) was almost constant across genome-wide significance levels of 1, 5 and 10%, while power to identifying true simulated QTL increased substantially as the significance level grew less stringent. Therefore, a more liberal (10%) genome-wide significance level could be used in genome-wide scan, corroborating the results of C. Laurie, S. Wang, L. A. Carlini-Garcia and Z-B. Zeng in MIM model (unpublished). Third, MTMIM model could have lower power than MIM model for QTL with effects on only a small subset of traits under analysis. However, as the subset of traits affected by a QTL increases, power in MTMIM model overpasses power in MIM model even when not all traits are affected by that QTL. Forth, MIM and MTMIM models revealed no qualitative differences in estimates of positions of QTL. However, there is a trend of smaller variation (measured in terms of standard error of means) in estimates of positions in MTMIM model for those QTL with pleiotropic effects on a larger subsets of traits under analysis. Fifth, the LOD-1.5 support interval produced confidence intervals for position of QTL with approximately 95% coverage in both MIM and MTMIM models. However, confidence interval for position of QTL was much wider in MIM than in MTMIM model. Sixth, estimates of effects of QTL in MTMIM model were in general closer to parameters for those QTL with effects on more traits. In MIM model effects of QTL were all over-estimated, demonstrating biases of selection due to genome-wide scan (Beavis, 1998). Overall, a qualitative comparison of results from

MIM and MTMIM models shows that effect estimates in the latter are less biased than in the former model. Lastly, LRT was shown to keep adequate type I error level when testing the null hypothesis of pleiotropic QTL against the alternative of close linked non-pleiotropic QTL in the bivariate analysis, while it delivered reasonable power when data were generated under the alternative.

Analysis of data from an experiment with *Drosophila* showed the potentials that our MTMIM model have to deliver complementary information for unveiling the genetic architecture of complex traits. Such as the potential to estimate effect of pleiotropic QTL, to test the hypothesis of pleiotropic QTL versus the alternative of close linked non-pleiotropic QTL, and to estimate the genetic covariance between traits.

Our results showed that the GEM-NR algorithm speeded up convergence of the likelihood function considerably when compared to the ECM algorithm, while still delivering stable parameter estimates.

Concluding, we showed empirically that the score-based threshold maintained type I error rate and false discovery rate within acceptable levels in MTMIM model. Furthermore, MTMIM model can be used as an extra tool to better extract information from multivariate data. Therefore, revealing with more details the genetic architecture of complex traits.

# 4

# Prediction of length of confidence interval for position of QTL in multiple trait analysis

In Chapter 1, we reviewed the current status of research on statistical methods for mapping multiple QTL in single and multiple complex traits within the maximum likelihood and Bayesian frameworks. In Chapter 2, we proposed a MTMIM model for QTL inference from inbred line crosses. We derived parameter estimators, extended the score-based method of ZOU *et al.* (2004) to estimate threshold in MTMIM model, proposed a forward selection to build a MTMIM model using the score-based threshold as the criterion to assess the significance of effects of QTL, and proposed a model optimization procedure. In Chapter 3, we evaluated our MTMIM model and score-based threshold method with regard to type I error, model fitting, and pleiotropy testing. We ended with analysis of data from an experiment with *Drosophila*. In this chapter, we derive analytical formulae for prediction of length of confidence interval for position of QTL and for prediction of shape of the LRT around the position of QTL in multiple trait analysis, under the assumption of infinitely many markers and large sample size. We end with some simulations to evaluate the length and coverage

of confidence interval for position of QTL in two linkage maps with distinct saturation of markers.

## 4.1   Introduction

Although maximum likelihood based methods provide both parameter estimates that are asymptotically unbiased and their variance-covariances, the variance of position QTL is often inaccurate (DARVASI and SOLLER, 1997) and hard to obtain because of the required second order derivatives of likelihood function. Therefore, practical application of maximum likelihood theory for construction of confidence interval for position of QTL is cumbersome. A practical "confidence interval" for position of QTL, the LOD-$d$ support interval (LANDER and BOTSTEIN, 1989), is defined as a continuous genomic region that includes the position of QTL and all positions on its left and right sides with LOD values larger than or equal to the LOD value at the position of QTL after subtraction of a positive constant $d$.

Under assumptions of infinitely many markers, large sample size, and normality distribution of phenotypic values of subjects, LANDER and BOTSTEIN (1989) showed that the LOD statistic follows a distribution proportional to a $\chi^2$ (LOD $\sim \frac{1}{2} \log_{10}(e) \chi^2_{df}$, where $e$ is the Euler's constant and $df$ stands for degrees of freedom). Therefore, if $df = 1$ (i.e., one parameter is been tested) the LOD-1 support interval of LANDER and BOTSTEIN (1989) leads to a predicted confidence interval with 96.8% coverage. However, it is been shown (VAN OOIJEN, 1992) that the same constant $d$ can lead to different levels of coverage depending on the effect of QTL and sample size. VAN OOIJEN (1992) pointed out that LOD-2 support interval would be necessary to achieve 95% coverage under their simulation settings. Moreover, MANGIN *et al.* (1994) showed evidences from simulation experiments that there are severe biasses in coverage of position of QTL with small effects, especially for saturate linkage maps, and the reason for these biasses is the poor convergence of the LRT statistic (or, equivalently the LOD statistic) towards a chi-squared distribution when effects

of QTL are small. In the light of their results, MANGIN *et al.* (1994) proposed a new statistic whose asymptotic distribution is independent of any nuisance parameter, especially the effects of QTL. However, practical use of the proposed statistic can be cumbersome because its computation burden and a correct threshold necessary for estimation of confidence interval must be found empirically through simulations. Nevertheless, their simulation experiments showed that confidence interval for position of QTL is unbiased for a wide range of effects of QTL (MANGIN *et al.*, 1994).

DARVASI and SOLLER (1997) defined a term "resolving power" as the 95% confidence interval (95CI) for position of QTL that one would obtain with infinitely many scored genetic markers. With many simulation experiments over a wide range of sample size ($n$) and allelic substitution effects ($\tau$), the authors empirically estimated a generalized equation (4.1) to predict resolving power of backcross ($m = 1$) and intercross $F_2$ ($m = 2$) populations.

$$95\mathrm{CI} = \frac{3000}{mn\tau^2} \tag{4.1}$$

DARVASI and SOLLER (1997) noticed that resolving power predicted under assumption of infinitely many markers was in general similar to 95CI for linkage maps with marker spaced 10 to 20 cM apart from each other.

The empirical results of DARVASI and SOLLER (1997) was latter analytically confirmed to hold good under large sample theory (VISSCHER and GODDARD, 2004). VISSCHER and GODDARD (2004) analytically derived equation (4.2), which is a generalization of equation (4.1). Equation (4.2) permits prediction of confidence interval for position of QTL with additive effect $a_t$ on trait $t$ in backcross ($x = 4$) and intercross $F_2$ ($x = 2$) populations for any given coverage level $100(1 - \alpha)$, where $0 \leq \alpha \leq 1$.

$$100\text{(1-}\alpha\text{)CI} = \frac{200xz_{\alpha/2}^2}{na_t^2} \tag{4.2}$$

where, for the standard normal distribution function with cumulative probability

distribution function $\Phi$, $z_{\alpha/2}$ satisfies $\Phi(z_{\alpha/2}) = 1 - \alpha/2$ (e.g., for $\alpha = 0.05$, $z_{\alpha/2} = 1.96$).

ROBERTS *et al.* (1999) showed evidence that standard error of position of QTL is a function of the expected LOD statistic. Furthermore, they empirically estimated equation (4.3) for prediction of standard error (SE) for position of QTL as a function of expected LOD (E[LOD]) in affected sib pair design.

$$SE = 7.0181[E(LOD)]^{-0.5881} \tag{4.3}$$

It is important to understand how this expression was generate before it can be applied. A wide range of genetic parameter settings were simulated mimicking real data. For each parameter setting, replicates were simulated and analyzed, producing LOD profile along the simulated chromosome. Then, an averaged LOD across replicates was computed for each position along the chromosome, leading to the averaged LOD profile along the chromosome. Next, the averaged LOD for positions around the *position* of the simulated QTL were used to fit a quadratic function of the form LOD$= b_0 + b_1 * position + b_2 * position^2$. Then, the square root of the reciprocal of the negative of the second order derivative of this function evaluated at the simulated position of QTL was taken as the measure of standard error for that particular parameter setting. Finally, equation (4.3) represents the fitted equation of SEs from all parameter settings regressed on their respective averaged LOD at the position of simulated QTL.

Care must me taken when estimating confidence interval for position of QTL thought equation (4.3) because estimated LOD are just poor estimates of the expected LOD (ROBERTS *et al.*, 1999). Nevertheless, ROBERTS *et al.* (1999)'s expression is still useful for a *priori* prediction of confidence interval for position of QTL with certain hypothesized effects, and so are the expressions derived by DARVASI and SOLLER (1997) and VISSCHER and GODDARD (2004).

ROBERTS *et al.* (1999) showed that estimated position of QTL with weak effect

may be many cM from its true position. ROBERTS *et al.* (1999) and other investigators (for instance, FLINT-GARCIA *et al.* (2003), CERVINO *et al.* (2005), and MACKAY and POWELL (2007)) argued for to use complementary approaches, such as linkage disequilibrium-based methods, to refine precision of estimated position of QTL.

While both DARVASI and SOLLER (1997) and VISSCHER and GODDARD (2004) have derived mathematical equations, empirically and analytically, respectively, to predict confidence interval for position of QTL with information from effect parameters and sample size, ROBERTS *et al.* (1999)'s empirical equation involves only the expected LOD. LOD summarizes all genetic information from any combination of experimental conditions, hence allowing an investigator to assess variation of position estimates across repeated experiments (ROBERTS *et al.*, 1999).

VISSCHER and GODDARD (2004) also studied shape of the LRT profile around the true position of QTL. Their results showed evidences for a non-quadratic form of LRT profile around the true position of QTL in saturate linkage maps. Therefore, contradicting asymptotic properties of the LRT statistic, from which one would expect a quadratic shape around the true QTL position, as it is been empirically shown in sparse maps (ROBERTS *et al.*, 1999).

In this chapter, we derive analytical formulae for prediction of length of confidence interval for position of QTL and for prediction of shape of LRT around the true position of QTL in multiple trait analysis, under assumptions of infinitely many markers and large sample size. Besides, we use simulation to evaluate length and coverage of confidence interval for position of QTL in linkage maps with distinct saturation of markers in single and multiple trait analyses.

## 4.2 Prediction of length of confidence interval for position of QTL

In this section, we derive analytical formulae for prediction of length of confidence interval for position of QTL in multiple trait analysis assuming infinitely many mark-

ers such that for any subject we know its genotype at any loci in the genome. With such saturate linkage map any QTL affecting a trait is assumed to coincide with a marker locus. For each trait $t$ ($t = 1, 2, \cdots, T$), we parameterize the genotypic values of genotypes $QQ$, $Qq$, $qq$ as $a_t$, $d_t$ and $-a_t$, respectively (FALCONER and MACKAY, 1996). We assume a BC population in which a QTL may assume genotype either $QQ$ or $Qq$, each with expected frequency of $\frac{1}{2}$. Similarly, any genetic marker may assume either genotype $MM$ or $Mm$, each with expected frequency of $\frac{1}{2}$. The expected genotypic variance of trait $t$ due to a QTL is $\frac{1}{4}(a_t - d_t)^2$, and the expected genotypic covariance of traits $t$ and $t'$ due to a QTL is $\frac{1}{4}(a_t - d_t)(a_{t'} - d_{t'})$ (see Appendix C).

The phenotypic measurement of trait $t$ on subject $i$, $y_{ti}$, ($i = 1, 2, \cdots, n$) can be written in the following linear model:

$$y_{ti} = \beta_t X_i + e_{ti} \tag{4.4}$$

where, $X_i$ takes value $\frac{1}{2}$ or $-\frac{1}{2}$, according to whether subject $i$ has genotype $QQ$ or $Qq$, respectively. The residual $e_{ti}$ is assumed to be independent and identically distributed according to a normal distribution with mean zero and variance $\sigma_{e_t}^2$. The vector of residuals $e_i = (e_{1i}, e_{2i}, \cdots, e_{Ti})'$ is assumed to follow a multivariate normal distribution with mean vector zero and variance-covariance $\boldsymbol{\Sigma}_e$. We define $\boldsymbol{X} = (X_1, X_2, \cdots, X_n)$, $\boldsymbol{\beta} = (\beta_1, \beta_2, \cdots, \beta_T)'$, $\boldsymbol{y}_i = (y_{1i}, y_{2i}, \cdots, y_{Ti})'$, $\boldsymbol{e}_i = (e_{1i}, e_{2i}, \cdots, e_{Ti})'$, $\boldsymbol{Y} = (\boldsymbol{y}_1, \boldsymbol{y}_2, \cdots, \boldsymbol{y}_n)$, and $\boldsymbol{E} = (\boldsymbol{e}_1, \boldsymbol{e}_2, \cdots, \boldsymbol{e}_n)$, and rewrite model (4.4) in matrix form as:

$$\boldsymbol{Y} = \boldsymbol{\beta} \boldsymbol{X} + \boldsymbol{E}$$

Under the assumption that residual error and genotypic effects are independent, the expected phenotypic variance of trait $t$ ($\sigma_{p_t}^2$) can be partitioned into expected genotypic ($\sigma_{g_t}^2$) and residual ($\sigma_{e_t}^2$) variances, i.e., $\sigma_{p_t}^2 = \sigma_{g_t}^2 + \sigma_{e_t}^2$. Similarly, the phenotypic covariance between traits $t$ and $t'$ can be partitioned in $\sigma_{p_{tt'}} = \sigma_{g_{tt'}} + \sigma_{e_{tt'}}$. The heritability of trait $t$ is defined as $h_t^2 = \sigma_{g_t}^2 / (\sigma_{g_t}^2 + \sigma_{e_t}^2)$. In matrix form, the phenotypic variance-covariance ($\boldsymbol{\Sigma}_p$) can be written as the sum of genotypic variance-

covariance ($\mathbf{\Sigma}_g$) and residual variance-covariance ($\mathbf{\Sigma}_e$).

$$
\begin{pmatrix}
\sigma_{p_1}^2 & \sigma_{p_{12}} & \cdots & \sigma_{p_{1T}} \\
\sigma_{p_{21}} & \sigma_{p_2}^2 & \cdots & \sigma_{p_{2T}} \\
\vdots & \vdots & \ddots & \vdots \\
\sigma_{p_{T1}} & \sigma_{p_{T2}} & \cdots & \sigma_{p_T}^2
\end{pmatrix}
=
\begin{pmatrix}
\sigma_{g_1}^2 & \sigma_{g_{12}} & \cdots & \sigma_{g_{1T}} \\
\sigma_{g_{21}} & \sigma_{g_2}^2 & \cdots & \sigma_{g_{2T}} \\
\vdots & \vdots & \ddots & \vdots \\
\sigma_{g_{T1}} & \sigma_{g_{T2}} & \cdots & \sigma_{g_T}^2
\end{pmatrix}
+
\begin{pmatrix}
\sigma_{e_1}^2 & \sigma_{e_{12}} & \cdots & \sigma_{e_{1T}} \\
\sigma_{e_{21}} & \sigma_{e_2}^2 & \cdots & \sigma_{e_{2T}} \\
\vdots & \vdots & \ddots & \vdots \\
\sigma_{e_{T1}} & \sigma_{e_{T2}} & \cdots & \sigma_{e_T}^2
\end{pmatrix}
$$

Assuming a large population we may observe four haplotypes $MQ$, $Mq$, $mQ$ and $mq$ for any marker-QTL pair. We assume that genetic marker and QTL are linked in *cis* phase (i.e., in $F_1$ generation alleles $M$ and $Q$ reside in the same chromosome) and have recombination frequency $r$. Under these assumptions, the four haplotypes have expected frequency $\frac{1-r}{2}$, $\frac{r}{2}$, $\frac{r}{2}$, $\frac{1-r}{2}$, respectively. For each trait $t$ and a given genetic marker we define the phenotypic mean of four haplotypes as $\bar{y}_{t_{MQ}}$, $\bar{y}_{t_{Mq}}$, $\bar{y}_{t_{mQ}}$, $\bar{y}_{t_{mq}}$. We collect means of traits for each haplotype into four vectors $\bar{\mathbf{Y}}_{MQ} = (\bar{y}_{1_{MQ}}, \bar{y}_{2_{MQ}}, \cdots, \bar{y}_{T_{MQ}})'$, $\bar{\mathbf{Y}}_{Mq} = (\bar{y}_{1_{Mq}}, \bar{y}_{2_{Mq}}, \cdots, \bar{y}_{T_{Mq}})'$, $\bar{\mathbf{Y}}_{mQ} = (\bar{y}_{1_{mQ}}, \bar{y}_{2_{mQ}}, \cdots, \bar{y}_{T_{mQ}})'$, and $\bar{\mathbf{Y}}_{mq} = (\bar{y}_{1_{mq}}, \bar{y}_{2_{mq}}, \cdots, \bar{y}_{T_{mq}})'$. We also define the phenotypic means of two genotypes $MM$ and $Mm$ as $\bar{y}_{t_M}$ and $\bar{y}_{t_m}$, respectively. We collect means of traits for genotypes $MM$ and $Mm$ into vectors $\bar{\mathbf{Y}}_M = (\bar{y}_{1_M}, \bar{y}_{2_M}, \cdots, \bar{y}_{T_M})'$ and $\bar{\mathbf{Y}}_m = (\bar{y}_{1_m}, \bar{y}_{2_m}, \cdots, \bar{y}_{T_m})'$, respectively. Likewise, we collect phenotypic means of traits for genotypes $QQ$ and $Qq$ into vectors $\bar{\mathbf{Y}}_Q = (\bar{y}_{1_Q}, \bar{y}_{2_Q}, \cdots, \bar{y}_{T_Q})'$ and $\bar{\mathbf{Y}}_q = (\bar{y}_{1_q}, \bar{y}_{2_q}, \cdots, \bar{y}_{T_q})'$, respectively. The phenotypic means of marker and QTL can be written is terms of haplotype means as follows:

$$
\begin{aligned}
\bar{\mathbf{Y}}_M &= (1-r)\bar{\mathbf{Y}}_{MQ} + r\bar{\mathbf{Y}}_{Mq} \\
\bar{\mathbf{Y}}_m &= (1-r)\bar{\mathbf{Y}}_{mq} + r\bar{\mathbf{Y}}_{mQ} \\
\bar{\mathbf{Y}}_Q &= (1-r)\bar{\mathbf{Y}}_{MQ} + r\bar{\mathbf{Y}}_{mQ} \\
\bar{\mathbf{Y}}_q &= (1-r)\bar{\mathbf{Y}}_{mq} + r\bar{\mathbf{Y}}_{Mq}
\end{aligned}
$$

In single trait analysis, under the assumption of infinitely many genetic markers the LRT statistic and F-statistic have been shown to produce very similar results when

testing for association between genetic marker and phenotypic variation of traits. The squared t-statistic is an alternative test, which is equivalent to F-statistic when only two treatment means are tested, i.e., only one degree of freedom in the denominator of F-statistic. Same relationship between squared t-statistic and F-statistic carries over to multiple trait analysis (RENCHER, 2002). Let $\hat{\mathbf{\Sigma}}_{(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)}$ be the estimated variance-covariance of difference between means $\bar{\boldsymbol{Y}}_M$ and $\bar{\boldsymbol{Y}}_m$, where $\bar{\boldsymbol{Y}}_M$ is the mean of $n_1$ subjects with genotype $MM$ and $\bar{\boldsymbol{Y}}_m$ is the mean of $n_2$ subjects with genotype $Mm$. The Hotelling's squared T-statistic (HOTELLING, 1951) to test for association between genetic marker and phenotypic variation of $T$ traits is:

$$T^2_{stat} = (\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)' \mathbf{\Sigma}^{-1}_{(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)} (\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)$$

We assume $\mathbf{\Sigma}_{(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)} \approx \frac{4}{n}\mathbf{\Sigma}_e$, which is reasonable under assumptions of infinitely many markers (small recombination between QTL and close linked genetic marker), small effects of QTL (LANDER and BOTSTEIN, 1989; LYNCH and WALSH, 1997; VISSCHER and GODDARD, 2004), and $n_1 = n_2 = \frac{n}{2}$.

With infinitely many markers the position of QTL is set to the position of the genetic marker with largest LRT, and confidence interval for position of QTL is constructed with the LOD-$d$ method (LANDER and BOTSTEIN, 1989) or bootstrap (VISSCHER et al., 1996). Our goal here is to derive an expression to predict length of confidence interval for position of QTL in multiple trait analysis. The key idea to build such confidence interval is based on the construction of a statistic (D-statistic, as in VISSCHER and GODDARD (2004)) as a function of the distance, $l$, between the marker with largest squared T-statistic and true position of QTL. Then, to apply inversion method (CASELLA and BERGER, 2002) to obtain length of confidence interval for position of QTL for any desired coverage level. We define the D-statistic as the difference between the squared T-statistic at a genetic marker positioned $l$ Morgans from the true position of the QTL and squared T-statistic at the true position of

QTL. The D-statistic is written as:

$$
\begin{aligned}
D(l) &= T^2_{stat}(l) - T^2_{stat}(0) \\
&= (\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)' \boldsymbol{\Sigma}^{-1}_{(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m) - (\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)' \boldsymbol{\Sigma}^{-1}_{(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)
\end{aligned}
$$

We assume $\boldsymbol{\Sigma}_{(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)} \approx \frac{4}{n}\boldsymbol{\Sigma}_e$. A QTL will be located off its true position if $D(l) > 0$. Therefore, knowing the distribution of $D(l)$ would allows us to analytically or numerically solve $P[D(l) > 0] = \frac{\alpha}{2}$ for $l$, and the solution $l$ provides the boundaries of $100(1 - \alpha)$ confidence interval for position of QTL. We show (Appendix C) that the D-statistic can be rewritten in terms of haplotype means as follows:

$$
D(l) = nr(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq})' \boldsymbol{\Sigma}^{-1}_e (\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})
$$

We also show (Appendix C) that for an additive genetic model (i.e., $d_t = 0$ for all $t \in \{1, 2, \cdots, T\}$), $E[D(l)] = -nr(1-r)\boldsymbol{a}'\boldsymbol{\Sigma}^{-1}_e\boldsymbol{a}$ and $Var[D(l)] = [nr(1-r)]^2 \frac{16}{n(1-e^{-4l})}\boldsymbol{a}'\boldsymbol{\Sigma}^{-1}_e\boldsymbol{a}$, where $\boldsymbol{a} = (a_1, a_2, \cdots, a_T)$. To simplify, we define $D^\star(l) = \frac{D(l)}{nr(1-r)}$, then $E[D^\star(l)] = -\boldsymbol{a}'\boldsymbol{\Sigma}^{-1}_e\boldsymbol{a}$ and $Var[D^\star(l)] = \frac{16}{n(1-e^{-4l})}\boldsymbol{a}'\boldsymbol{\Sigma}^{-1}_e\boldsymbol{a}$. Therefore, $P[D(l) > 0] = P[D^\star(l) > 0]$.

Under large sample size assumption, $Z = \frac{D^\star(l) - E[D^\star(l)]}{\sqrt{Var[D^\star(l)]}}$ is normally distributed with mean zero and unit variance, and $P(D^\star(l) > 0) \approx P(Z > -\frac{E[D^\star(l)]}{\sqrt{Var[D^\star(l)]}})$. Therefore, the boundaries of the $100(1 - \alpha)$CI for position of QTL can be obtained as follows:

$$
\begin{aligned}
P(Z > -\frac{E[D^\star(l)]}{\sqrt{Var[D^\star(l)]}}) &= \frac{\alpha}{2} \\
-\frac{E[D^\star(l)]}{\sqrt{Var[D^\star(l)]}} &= z_{\alpha/2} \\
\frac{\boldsymbol{a}'\boldsymbol{\Sigma}^{-1}_e\boldsymbol{a}}{[\frac{16}{n(1-e^{-4l})}\boldsymbol{a}'\boldsymbol{\Sigma}^{-1}_e\boldsymbol{a}]^{\frac{1}{2}}} &= z_{\alpha/2}
\end{aligned}
\tag{4.5}
$$

Solving equation (4.5) for $l$ gives the boundaries of the confidence interval for

position of QTL:

$$l = -\frac{1}{4} \log_e(1 - z^2_{\alpha/2} \frac{16}{n\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}}) \tag{4.6}$$

Using the Taylor's approximation of $\log_e(1 - x) \approx -x$, for small $x$, equation (4.6) can be simplified to:

$$l = z^2_{\alpha/2} \frac{4}{n\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}}$$

The length of the $100(1 - \alpha)$CI, in cM, is:

$$(1 - \alpha)\text{CI} = 100l$$

$$= 400\frac{z^2_{\alpha/2}}{n\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}}$$

Because the approximation of length of confidence interval given by equation (4.5) does not account for the length of the chromosome in which the QTL is located, the approximation can be biased upwards, or even worse, it can produce confidence interval wider than the whole length of the chromosome. The upwards bias appears because the left hand side of equation (4.5) evaluated with $l$ at the ends of the chromosome can return values that are smaller than $z_{\alpha/2}$. Here, we propose a method for prediction of length of confidence interval for position of QTL that accounts for length of chromosome. The idea is to use a truncated normal distribution rather than a normal distribution. We take the upper limit of the truncated normal distribution as the minimum (min) of two values returned by the left hand side of equation (4.5) when evaluated at the ends of the chromosome. Let $l_1$ and $l_2$ denote the distance from the QTL to the left and right ends of chromosome respectively. Then, the upper limit ($u$) is:

$$u = \min_{l \in \{l_1, l_2\}} \left\{ \frac{\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}}{[\frac{16}{n(1-e^{-4l})}\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}]^{\frac{1}{2}}} \right\} \tag{4.7}$$

We take the lower limit of the truncated normal distribution as the negative

of the upper limit. Since the position of QTL in the chromosome has little effect in the truncation limits we can assume the position of QTL in the middle of the chromosome. The boundaries of the adjusted confidence interval for position of QTL can be obtained from:

$$\frac{\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}}{[\frac{16}{n(1-e^{-4l})}\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}]^{\frac{1}{2}}} = z_{u_{\alpha/2}} \tag{4.8}$$

where, $z_{u_{\alpha/2}}$ is the quantile of a truncated normal distribution with mean zero and unit variance and truncation points $-u$ and $u$. Then, using similar arguments as previously, the adjusted length of the $100(1\text{-}\alpha)$CI, in cM, is:

$$(1-\alpha)\mathrm{CI}_a = 400\frac{z_{u_{\alpha/2}}^2}{n\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}}$$

Hereafter, we denoted the method that uses truncated normal distribution as adjusted method and the other method as unadjusted. Although the adjusted method is very naive we empirically show (next sections) that it greatly improves the estimates of length of confidence interval for position of QTL when compared to the unadjusted method. Moreover, the adjusted method always guarantees length of confidence interval for position of QTL shorter than the whole length of the chromosome.

Let $f_t$ denote the fraction of genotypic variation of trait $t$ explained by a QTL with additive effect only. Specifically, $f_t = \frac{a_t^2}{4\sigma_{g_t}^2}$. Then, after some algebra, we have $4\sigma_{e_t}^2/a_t^2 = 1/f_t \times (1-h_t^2)/h_t^2$. Let the residual correlation between two traits be $\rho = \frac{\sigma_{e_{12}}}{\sqrt{\sigma_{e_1}^2 \sigma_{e_2}^2}}$. Therefore, the $(1-\alpha)$CI for position of QTL in bivariate analysis is:

$$\begin{aligned}(1-\alpha)\mathrm{CI} &= 100\frac{z_{u_{\alpha/2}}^2(1-\rho^2)}{\frac{n}{4}(\frac{a_1^2}{\sigma_{e_1}^2} - sign(a_1a_2)2\rho\frac{a_1a_2}{\sigma_{e_1}\sigma_{e_2}} + \frac{a_2^2}{\sigma_{e_2}^2})} \\ &= 100\frac{z_{u_{\alpha/2}}^2(1-\rho^2)}{n(\frac{f_1h_1^2}{1-h_1^2} - sign(a_1a_2)2\rho\frac{\sqrt{f_1h_1^2}}{\sqrt{1-h_1^2}}\frac{\sqrt{f_2h_2^2}}{\sqrt{1-h_2^2}} + \frac{f_2h_2^2}{1-h_2^2})}\end{aligned}$$

where, $sign(x_0) = 1$ if $x_0 \geq 0$ and $sign(x_0) = -1$ otherwise.

## 4.3   Prediction of shape of LRT around the position of QTL

In this section, we derive analytical formulae to predict shape of the LRT around the position of QTL in multiple trait analysis assuming infinitely many markers. Under the assumption of infinitely many markers, the likelihood function of data under model (4.4) does not involve mixture of distributions, which greatly simplifies our derivations. The likelihood function of model (4.4) at a genetic marker located $l$ Morgans apart from the position of a QTL and assuming an additive genetic model ($d_t = 0$ for all $t \in \{1, 2, \cdots, T\}$), is:

$$L(\boldsymbol{a}, \boldsymbol{\Sigma}_e | \boldsymbol{Y}, \boldsymbol{X}, l) = \prod_{i=1}^{n} (2\pi)^{-\frac{T}{2}} |\boldsymbol{\Sigma}_e|^{-\frac{1}{2}} e^{-\frac{1}{2}(\boldsymbol{y}_i - \boldsymbol{a}X_i)' \boldsymbol{\Sigma}_e^{-1}(\boldsymbol{y}_i - \boldsymbol{a}X_i)} \tag{4.9}$$

The natural logarithm of the likelihood function (4.9), $\ell(l) = \log_e L(\boldsymbol{a}, \boldsymbol{\Sigma}_e | \boldsymbol{Y}, \boldsymbol{X}, l)$, is:

$$\ell(l) = -\frac{nT}{2} \log(2\pi) - \frac{n}{2} \log(|\boldsymbol{\Sigma}_e|) - \sum_{i=1}^{n} \frac{1}{2}(\boldsymbol{y}_i - \boldsymbol{a}X_i)' \boldsymbol{\Sigma}_e^{-1}(\boldsymbol{y}_i - \boldsymbol{a}X_i) \tag{4.10}$$

The MLE of $\boldsymbol{\Sigma}_e$ and $\boldsymbol{a}$ at a given genetic marker $l$ Morgans apart from the position of a QTL are:

$$\hat{\boldsymbol{\Sigma}}_e(l) = \frac{1}{n}(\boldsymbol{Y} - \hat{\boldsymbol{a}}\boldsymbol{X})(\boldsymbol{Y} - \hat{\boldsymbol{a}}\boldsymbol{X})'$$

$$\hat{\boldsymbol{a}}(l) = \boldsymbol{Y}\boldsymbol{X}'(\boldsymbol{X}\boldsymbol{X}')^{-1}$$

$$= \bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m$$

The last equality is true under assumption of equal number of subjects with genotypes

$MM$ and $Mm$ (i.e., $n_1 = n_2 = n/2$). The logarithm likelihood function (4.10) evaluated at the MLE is:

$$\hat{\ell}(l) = -\frac{nT}{2}\log(2\pi) - \frac{n}{2}\log(|\hat{\mathbf{\Sigma}}_e(l)|) - \frac{n}{2}$$

The shape of LRT around the position of a QTL can be assessed by means of the expectation $\mathrm{E}[\mathrm{LRT}(l)] = -2\mathrm{E}[\hat{\ell}(l) - \hat{\ell}(0)]$. We show (Appendix C) that shape of LRT around the position of a QTL can be predicted by:

$$\mathrm{E}[\mathrm{LRT}(l)] \approx \frac{n}{4}(1 - e^{-4l})\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a} \tag{4.11}$$

For the bivariate analysis we can rewrite equation (4.11) as:

$$\begin{aligned}
\mathrm{E}[\mathrm{LRT}(l)] &= \frac{n(1 - e^{-4l})}{4(1 - \rho^2)}\left(\frac{a_1^2}{\sigma_{e_1}^2} - sign(a_1 a_2)2\rho\frac{a_1 a_2}{\sigma_{e_1}\sigma_{e_2}} + \frac{a_2^2}{\sigma_{e_2}^2}\right) \\
&= \frac{n(1 - e^{-4l})}{(1 - \rho^2)}\left(\frac{f_1 h_1^2}{1 - h_1^2} - sign(a_1 a_2)2\rho\frac{\sqrt{f_1 h_1^2}\sqrt{f_2 h_2^2}}{\sqrt{1 - h_1^2}\sqrt{1 - h_2^2}} + \frac{f_2 h_2^2}{1 - h_2^2}\right)
\end{aligned}$$

Therefore, rather than knowing the QTL effect and environmental variation in the bivariate analysis, if we know the heritability and the contribution of the QTL to the genotypic variance we can estimate the shape of the LRT around the QTL.

## 4.4  Simulations

In this section, we use simulation to assess empirical length of confidence interval for position of QTL constructed with the LOD-$d$ method. In our simulations, we varied heritability level of QTL (5, 10 and 15% of phenotypic variation), number of traits (one, two and three traits), and marker distance (2 and 10 cM, 2-cM map and 10-cM map, respectively). In the 2-cM map, we simulated one thousand replicates with parameters shown in Table 4.1. Each replicate consisting of one chromosome with 41 markers evenly spaced at 2 cM apart from each other, three traits and sample

size of 300. In the genome-wide scan we tested the presence of a putative QTL at each genetic marker. In the 10-cM map, we simulated one thousand replicates with parameters shown in Table 4.1. Each replicate consisting of one chromosome with 9 markers evenly spaced at 10-cM apart from each other, 300 subjects and three traits. In the genome-wide scan we tested the presence of a putative QTL at each 1-cM in the genome. In both experiments, 2- and 10-cM maps, we used the score-based threshold to assess the genome-wide significance of effects of putative QTL. The score statistic was resampled 800 times at each tested position in the genome.

Table 4.1: Parameters used in simulations to evaluate length of confidence interval for position of QTL.

| Traits | Effects of the QTL | | | $\mathbf{\Sigma}_e$[$\$] | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | $h^2=.05$ | $h^2=0.10$ | $h^2=0.15$ | T1 | T2 | T3 |
| T1 | 0.46 | 0.67 | 0.84 | 1 | 0.2 | 0 |
| T2 | 0.46 | 0.67 | 0.84 | 0.2 | 1 | -0.2 |
| T3 | 0.46 | 0.67 | 0.84 | 0 | -0.2 | 1 |
| Position[b] | 34 | 34 | 34 | – | – | – |

[$] Residual variance-covariance matrix.
[b] Position, in cM, of QTL from the leftmost marker in the chromosome.

## 4.5   Results

We evaluated the unadjusted (equation (4.5)) and adjusted (equation (4.8)) methods to predict length of confidence interval for position of QTL. In both methods, we varied number of traits, heritability levels of QTL and sample sizes (Table 4.1). In the adjusted method, the lower and upper limits of the truncated normal distribution were obtained by equation (4.7) assuming that a QTL is located in the middle of the simulated 80-cM long chromosome ($l_1 = l_2 = 40/100$ Morgans) and parameters in Table 4.1. Our results show a remarkable difference between unadjusted (Figure 4.1A) and adjusted (Figure 4.1B) methods for predicting length of confidence interval for position of QTL with heritability 5%, sample size 300 and single trait analysis. While length of confidence interval predicted by the unadjusted method was approximately 90.49 cM (wider than the whole length of the chromosome–80 cM), the the length of confidence interval predicted by the adjusted method was approximately 23.26 cM. Table 4.2 shows empirical (from simulations) and analytical (unadjusted and adjusted methods) predicted length of confidence interval for position of QTL. The results clearly demonstrate that while the empirical and adjusted estimates of length of confidence intervals are very similar for all levels of heritability and number of traits, the empirical and unadjusted estimates can be remarkably different in single trait analysis when a QTL has low heritability.

Graphical visualization (Figure 4.1B) facilitates understanding the role of each factor involved in predicting the length of confidence interval for position of QTL. First, there is a inverse relationship between percentage of phenotypic variation explained by a QTL and length of confidence interval, for any population size and number of traits analyzed jointly. Second, length of confidence interval decreases when population size increases, for any heritability level of QTL and any number of traits analyzed jointly. And lastly, if a QTL has pleiotropic effects, multiple trait analysis delivers more precise estimates of position of QTL.

Figure 4.1: Prediction of length of 95% confidence interval for position of QTL with infinitely many markers. Predictions are displayed as function of heritability of QTL, sample size (n) and number of traits (T) analyzed jointly. (A) and (B) show predicted length of confidence interval delivered by unadjusted (equation 4.5) and adjusted methods (equation 4.8), respectively.

We also evaluated equation (4.11) to assess the shape of LRT around position of QTL. We varied heritability levels (5 and 15%) and number of traits (one, two and three). Two factors, number of traits and heritability of QTL, have important role in shaping the LRT (Figures 4.2A and 4.2C). As the testing position moves away from the position of the QTL, the LRT decays more rapidly when more traits are analyzed jointly and/or when the QTL has higher heritability. Therefore, multiple trait analysis and/or stronger effects can led to narrower confidence interval for position of QTL, provided that the QTL has pleiotropic effects. This result corroborates the results displayed in Figure 4.1. Analytical predictions of shape of LRT around the position of QTL agree very well with the averaged LRT over 1000 replicates (Figures 4.2B and

Table 4.2: Analytical and empirical length of 95% confidence interval for position of QTL in single and multiple trait analyses. Three heritability ($h^2$) were investigated.

| Analysis | Length | $h^2$=5% | $h^2$=10% | $h^2$=15% |
|---|---|---|---|---|
| | empirical[*] | 26.3 (1.06) | 13.3 (0.50) | 7.9 (0.29) |
| Single (T1) | adjusted[$] | 23.3 | 13.9 | 8.5 |
| | unadjusted[$] | 90.5 | 15.5 | 8.6 |
| | empirical | 18.7 (0.72) | 7.1 (0.26) | 4.1 (0.14) |
| Multiple (T1,T2) | adjusted | 17.1 | 8.0 | 4.8 |
| | unadjusted | 21.9 | 8.1 | 4.8 |
| | empirical | 9.9 (0.36) | 4.1 (0.14) | ** |
| Multiple (T1,T2,T3) | adjusted | 9.3 | 4.1 | 2.5 |
| | unadjusted | 9.5 | 4.1 | 2.5 |

[*] Average of length of empirical 95% confidence intervals built with the LOD-$d$ method. The value of $d$ in LOD-$d$ was fine tuned until LOD-$d$ support interval of 95% of all replicates that showed significant QTL effect included the simulated QTL. Then, length of LOD-$d$ support interval for those replicates with LOD-$d$ support interval including the simulated QTL were averaged. Standard error of means are in parentheses.

[$] Adjusted and unadjusted analytical approximations of length of 95% confidence interval.

[**] Even very small values of $d$ in LOD-$d$ led to empirical confidence interval with coverage larger than 98%. Therefore, we were unable to verify the length of empirical 95% confidence interval.

4.2D). Therefore, our analytical predictions might be useful to predict shape of LRT around a hypothesized QTL.

Equation (4.11) can be decomposed into two non-centrality parameters of two chi-squared probability distribution functions. When the LRT is performed at the position of a QTL, the LRT follows a chi-squared distribution with non-centrality parameter $\frac{n}{4}\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}$ and degrees of freedom equal to difference in number of parameters in the full and reduced models. We demonstrate that the LRT, when evaluated at a genomic position $l$ Morgans (Haldane distance) from the position of QTL, follows a chi-squared distribution with non-centrality parameter $\frac{n}{4}e^{-4l}\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}$. In Figure 4.3, we show the empirical distribution of LRT at various positions in a simulated chromosome (simulations setting of 2-cM map). These results regard the multiple trait analysis of traits T1 and T2 with a pleiotropic QTL positioned 34 cM from the leftmost marker in the chromosome and with heritability of 5% (see Table 4.1).

Figure 4.2: Analytical (A and C from equation (4.11)) and averaged LRT over 1000 replicates (B and D from simulations) prediction of shape of LRT around the position of QTL. Prediction of shape of LRT is displayed as function of heritability level of QTL (A and B: $h^2$=5%, C and D: $h^2$=15%), number of traits (T), and distance ($l$) from testing position and position of QTL. A linkage map with genetic markers 2 cM apart from each other and sample size n=300 were assumed in both simulation and analytical predictions.

Figure 4.3: Empirical distribution of LRT at three positions in a simulated chromosome. These results regard the multiple trait analysis of traits T1 and T2 with a pleiotropic QTL of heritability 5%. From top to bottom, first row shows the histogram (left panel) of LRT values exactly at the position of QTL, overlayed by a chi-squared probability distribution function with two degrees of freedom and non-centrality parameter $\xi = \frac{n}{4}e^{-4l}\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}$ ($\chi_2^2(\xi)$), where $l$ is distance of the testing position from the position of QTL, $n$ is the sample size, $\boldsymbol{a}$ is the vector of effects of QTL and $\boldsymbol{\Sigma}_e$ is the residual variance-covariance matrix. The right panel shows a quantile-quantile plot of a $\chi_2^2(\xi)$ distribution versus the LRT. Second row shows results of LRT 10 cM away from the position of QTL, which leads to $l = 10/100$. Third row shows results of LRT 20 cM away from the position of QTL, which leads to $l = 20/100$. A linkage map with genetic markers 2 cM apart from each other and sample size n=300 were assumed in simulations. A total of 1000 replicates were analyzed.

Results of simulations of a QTL explaining 5% of phenotypic variation of each trait in 2-cM and 10 cM maps are summarized in Table 4.3. There are three factors (distance between markers, number of traits analyzed jointly, and level $d$ in LOD-$d$ method) that need to be discussed from this experiment. First, for any level $d$ in LOD-$d$ and number of traits, length of confidence interval for position of QTL in 10-cM map is wider than in 2-cM map. As an immediate consequence, the confidence interval in 10-cM map has higher coverage. Figures 4.2B and 4.4 show clearly that shape of LRT around the position of QTL in 10-cM map is flatter than in 2-cM map. Flatness and spikiness of LRT around position of QTL impose wider and shorter confidence intervals in 10- and 2-cM maps, respectively.

Our results provide empirical evidences that in order to keep approximately 95% confidence interval in both 10-cM and 2-cM maps, $d$ levels in LOD-$d$ should be 1 and 1.5, respectively, regardless of number of traits analyzed jointly. Two distinct levels of $d$ for different maps is mainly due to differences in shape of LRT around the position of QTL.

Our results provide some evidences that a region in the genome that gives approximately 95% of chance to include the QTL is not much shorter in 2-cM map as compared to 10-cM map if single trait analysis is performed (in our simulations, the length of confidence interval is approximately 26 cM in both 2- and 10-cM maps). On the other hand, there seems to exist a minor gain in precision when we moved from 10- to 2-cM maps if multiple trait analysis is performed and the putative QTL has pleiotropic effects. For instance, in our simulations, multiple trait analysis of three traits with a small pleiotropic QTL led to length of confidence interval for position of QTL of 14.0 and 10.5 cM, in 10- and 2-cM maps, respectively.

Within both 2- and 10-cM maps, there is major reduction in length of confidence interval for position of pleiotropic QTL in multiple trait analysis. For instance, the LOD-1.5 support intervals in 2-cM map were 26.3, 17.5 and 10.5 cM, according to whether one, two and three traits were analyzed jointly, respectively.

Table 4.3: Coverage (%) and length (cM) of LOD-$d$ support interval for a QTL explaining 5% of phenotypic variation of traits. Estimated position of QTL are also shown. Single and multiple trait analyses were performed with 5% genome-wide significance level. Standard errors are in parentheses.

| Map | Summary | Traits | LOD-$d$ level | | | | |
|-----|---------|--------|------|------|------|------|------|
| | | | 0.50 | 0.75 | 1 | 1.5 | 2 |
| 2-cM | Coverage | 1 | 70.2 | 80.0 | 86.5 | 95.0 | 97.9 |
| | | 1,2 | 77.4 | 83.4 | 87.2 | 94.4 | 97.6 |
| | | 1,2,3 | 87.9 | 89.9 | 92.1 | 95.9 | 98.8 |
| | Length | 1 | 8.0 (0.35) | 11.5 (0.49) | 15.8 (0.65) | 26.3 (1.06) | 38.1 (1.49) |
| | | 1,2 | 6.9 (0.28) | 9.0 (0.35) | 11.4 (0.45) | 17.5 (0.68) | 24.5 (0.95) |
| | | 1,2,3 | 5.6 (0.20) | 6.6 (0.24) | 7.7 (0.28) | 10.5 (0.38) | 13.7 (0.49) |
| | Position | 1 | 36.0 (0.10) | 36.1 (0.14) | 36.0 (0.18) | 35.5 (0.23) | 35.4 (0.24) |
| | | 1,2 | 35.9 (0.08) | 35.9 (0.10) | 36.0 (0.11) | 35.9 (0.15) | 35.8 (0.17) |
| | | 1,2,3 | 36.0 (0.05) | 36.0 (0.05) | 36.0 (0.06) | 36.0 (0.08) | 36.0 (0.09) |
| 10-cM | Coverage | 1 | 76.7 | 87.6 | 94.3 | 99.2 | 99.7 |
| | | 1,2 | 81.8 | 89.9 | 93.7 | 98.5 | 99.8 |
| | | 1,2,3 | 83.0 | 89.9 | 94.0 | 98.5 | 99.9 |
| | Length | 1 | 15.9 (0.24) | 20.8 (0.32) | 26.0 (0.38) | 36.8 (0.55) | 48.4 (0.67) |
| | | 1,2 | 12.8 (0.15) | 16.7 (0.21) | 20.2 (0.29) | 27.0 (0.38) | 34.5 (0.50) |
| | | 1,2,3 | 9.3 (0.09) | 11.7 (0.12) | 14.0 (0.16) | 18.4 (0.23) | 22.4 (0.28) |
| | Position | 1 | 35.4 (0.18) | 35.5 (0.21) | 35.4 (0.23) | 35.2 (0.26) | 35.0 (0.26) |
| | | 1,2 | 35.5 (0.15) | 35.6 (0.16) | 35.6 (0.17) | 35.5 (0.18) | 35.3 (0.20) |
| | | 1,2,3 | 35.6 (0.09) | 35.7 (0.10) | 35.7 (0.11) | 35.8 (0.13) | 35.8 (0.13) |

# 4.6   Concluding remarks

We proposed analytical formulae that allow prediction of length of confidence interval for position of QTL and prediction of shape of the LRT around the position of QTL with infinitely many markers and multiple trait analysis using large sample theory. Our results generalize the results of VISSCHER and GODDARD (2004) and they can be used to predict the length of confidence interval for position of QTL with a hypothesized effect on multiple trait, for any given coverage probability. Our analytical formulae can also be used to predict shape of LRT around the position of QTL. Furthermore, we have proposed an alternative method for predicting the length of confidence interval for position of QTL, the adjusted method. The adjusted method accounts for the length of the chromosome in which the QTL is located and we showed empirically that it can deliver more accurately predictions than the method with no adjustments, especially for QTL of low heritability. Our simulation results showed that for sample size of 300 and QTL with heritability levels of 5, 10 and 15%, there are good agreement between length of confidence intervals empirically estimated and analytically predicted with the adjusted method.

As in DARVASI and SOLLER (1997), we have shown that because of inverse relationship between length of confidence interval for position of QTL (resolving power) and sample size, major steps towards increasing resolution of QTL mapping is possible by increasing sample size. Moreover, we demonstrate that shorter confidence intervals can also be obtained by multiple trait analysis if the putative QTL has pleiotropic effects. However, as pointed out by DARVASI and SOLLER (1997), it is worth mentioning that for putative QTL that has power less than one, which means that if several experiments are repeated the QTL would not be identified in all of them, effects of QTL are overestimated. Therefore, estimated length of confidence interval for position of QTL in real data analysis is usually shorter than predicted.

Our simulations results showed that length of 95% confidence interval in 2- and 10-cM maps were approximately the same regardless of whether single or multiple trait

analyses were carried out. This finding corroborates results of Roberts *et al.* (1999) and Darvasi and Soller (1997) that decreasing marker distance below resolution power of QTL without increasing sample size does not improve estimation of position of QTL. Despite existence of empirical evidences showing little reduction in length of confidence interval for position of QTL when number of genetic markers is larger, use of highly saturate linkage maps, such as single nucleotide polymorphism (SNP) maps, is certainly advantageous because much more information is conveyed with such saturated maps. Therefore, combining such saturated linkage maps, large sample size, well-designed experiments, linkage analysis, linkage disequilibrium analysis, and information on gene pathways in the region surrogating the QTL might provide a way of narrowing down the search for candidate genes regulating traits of interest (Mackay, 2001; Cervino *et al.*, 2005).

In the 10-cM map, we provided empirical evidences that shape of LRT is quadratic around the position of QTL position in multiple trait analysis. Therefore, corroborating results of single trait analysis shown by Roberts *et al.* (1999).

In the 2-cM map, we provided analytical and empirical evidences that shape of LRT shape is not quadratic in multiple trait analysis. Therefore, corroborating results of single trait analysis shown by Visscher and Goddard (2004). Moreover, we demonstrated that in the presence of QTL, the LRT in the neighborhood of the QTL follows a chi-squared distribution with non-centrality parameter that depends on the distance between the testing position and the position of the QTL, sample size, effects of QTL, and structure of residual variance-covariance matrix.

We reinforce findings of Visscher and Goddard (2004) that asymptotic theory of LRT may not hold good in highly saturate linkage maps because the number of recombinants has a critical role in convergence. The formula of $D^{\star}(l)$ let clear that as the distance between marker decreases ($l$) goes to zero, the variance of $D^{\star}(l)$ blows up unless the number of recombinants ($nr$) is large. However, as distance between markers decreases, the number of recombinants between any two markers will reduce unless sample size grows very large. Therefore, since $D^{\star}(l)$ and LRT are equivalent

in convergence, asymptotic inference based on LRT distribution is cumbersome and care should be taken when drawing conclusions from LRT asymptotic theory in highly saturate linkage maps.

Figure 4.4: Empirical (from simulations) prediction of shape of LRT around the position of a pleiotropic QTL in 10-cM map. The prediction of shape is displayed as function of number of traits (T) and distance ($l$) from testing position and position of QTL. The pleiotropic QTL explained 5% of the phenotypic variation of each trait. One thousand replicates with sample size n=300 were used.

# Bibliography

AITKIN, M. and I. AITKIN, 1996 A hybrid EM/Gauss-Newton algorithm for maximum likelihood in mixture distributions. Statistics and Computing **6**: 127–130.

AKAIKE, H., 1969 A method of statistical identification of discrete time parameter linear systems. Annals of the Institute of Statistical Mathematics **21**: 243–247.

BALDING, D., M. BISHOP, and C. CANNINGS, editors, 2007 *Handbook of Statistical Genetics*, volume 1. Wiley, third edition.

BANERJEE, S., B. S. YANDELL, and N. JUN YI, 2008 Bayesian quantitative trait loci mapping for multiple traits. Genetics **179**: 2275–2289.

BEAVIS, W. D., 1998 QTL analyses: power, precision, and accuracy, pp. 145–162 in *Molecular Dissection of Complex Traits*, edited by PATERSON, A. H., CRC Press, New York.

BENJAMINI, Y. and Y. HOCHBERG, 1995 Controlling the false discovery rate: A practical and powerful approach to multiple testing. Journal of the Royal Statistical Society. Series B (Methodological) **57**: 289–300.

BROMAN, K. W. and T. P. SPEED, 2002 A model selection approach for the identification of quantitative trait loci in experimental crosses. Journal of the Royal Statistical Society. Series B (Statistical Methodology) **64**: 641–656.

CASELLA, G. and R. L. BERGER, 2002 *Statistical Inference*. Duxbury Press, second edition.

CERVINO, A. C., G. LI, S. EDWARDS, J. ZHU, C. LAURIE, G. TOKIWA, P. Y. LUM, S. WANG, L. W. CASTELLINI, A. J. LUSIS, S. CARLSON, A. B. SACHS, and E. E. SCHADT, 2005 Integrating qtl and high-density snp analyses in mice to identify insig2 as a susceptibility gene for plasma cholesterol levels. Genomics **86**: 505–517.

CHURCHILL, G. A. and R. W. DOERGE, 1994 Empirical threshold values for quantitative trait mapping. Genetics **138**: 963–971.

COX, D. R. and D. V. HINKLEY, 1974 *Theoretical Statistics*. Chapman and Hall, London.

DARVASI, A. and M. SOLLER, 1997 A simple method to calculate resolving power and confidence interval of qtl map location. Behavior Genetics **27**: 125–132.

DAVIES, R., 1977 Hypothesis testing when nuisance parameter is present only under the alternative. Biometrika **64**: 247–254.

DAVIES, R. B., 1987 Hypothesis testing when a nuisance parameter is present only under the alternative. Biometrika **74**: 33–43.

DEMPSTER, A. P., N. LAIRD, and D. RUBIN, 1977 Maximum likelihood from incomplete data via the em algorithm. Journal of the Royal Statistical Society. Series B (Methodological) **39**: 1–38.

DOERGE, R. W. and G. A. CHURCHILL, 1996 Permutation tests for multiple loci affecting a quantitative character. Genetics **142**: 285–294.

DWYER, P. S., 1967 Some applications of matrix derivatives in multivariate analysis. Journal of the American Statistical Association **62**: 607–625.

DWYER, P. S. and M. S. MACPHAIL, 1948 Symbolic matrix derivatives. The Annals of Mathematical Statistics **19**: 517–534.

EAVES, L. J., M. C. NEALE, and H. MAES, 1996 Multivariate multipoint linkage analysis trait quantitative trait loci. Behavior Genetics **26**: 519–525.

EFRON, B. and R. J. TIBSHIRANI, 1993 *An introcuction to the Bootstrap*. Chapman and Hall/CRC.

FALCONER, D. and T. F. MACKAY, 1996 *Introduction to quantitative genetics*. Longman, fourth edition.

FALCONER, D. S., 1952 The problem of environment and selection. Amer. Nat. **86**: 293–298.

FLINT-GARCIA, S. A., J. M. THORNSBERRY, and E. S. B. IV, 2003 Structure of linkage disequilibrium in plants. Annual Review of Plant Biology **54**: 357–374.

GARCIA, A. A. F., S. WANG, A. E. MELCHINGER, and Z.-B. ZENG, 2008 Quantitative trait loci mapping and the genetic basis of heterosis in maize and rice. Genetics **180**: 1707–1724.

GEMAN, S. and D. GEMAN, 1984 Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. IEEE Transactions on Pattern Analysis and Machine Intelligence **6**: 721–741.

GODSILL, S. J., 2001 On the relationship between markov chain monte carlo methods for model uncertainty. Journal of Computational and Graphical Statistics **10**: 1–19.

GREEN, P. J., 1995 Reversible jump markov chain monte carlo computation and bayesian model determination. Biometrika **82**: 711–732.

HALEY, C. S. and S. A. KNOTT, 1992 A simple regression method for mapping quantitative trait loci in line crosses using flanking markers. Heredity **69**: 315–324.

HASTINGS, W. K., 1970 Monte carlo sampling methods using markov chains and their applications. Biometrika **57**: 97–109.

HOESCHELE, I., P. UIMARI, F. E. GRIGNOLA, Q. ZHANG, and K. M. GAGE, 1997 Advances in statistical methods to map quantitative trait loci in outbred populations. Genetics **147**: 1445–1457.

HOESCHELE, I. and P. VANRADEN, 1993 Bayesian analysis of linkage between genetic markers and quantitative trait loci. i. prior knowledge. Theor Appl Genet **85**: 953–960.

HOSCHELE, I., 2007 Mapping outbred quantiative trait loci outbred pedigrees, pp. 623–677 in *Handbook of Statistical Genetics*, edited by BALDING, D., M. BISHOP, and C. CANNINGS, Wiley.

HOTELLING, H., 1951 A generalized t test and measure of multivariate dispersion. In *Proc. Second Berkeley Symp. on Math. Statist. and Prob.*, pp. 23–41, Univ. of Calif. Press.

JANSEN, R. C., 1993 Interval mapping of multiple quantitative trait loci. Genetics **135**: 205–211.

JIANG, C. and Z. B. ZENG, 1995 Multiple trait analysis of genetic mapping for quantitative trait loci. Genetics **140**: 1111–1127.

JIANG, C. and Z. B. ZENG, 1997 Mapping quantitative trait loci with dominant and missing markers in various crosses from two inbred lines. Genetica **101**: 47–58.

KAO, C. H. and Z. B. ZENG, 1997 General formulas for obtaining the mles and the asymptotic variance-covariance matrix in mapping quantitative trait loci when using the em algorithm. Biometrics **53**: 653–665.

KAO, C.-H. and Z.-B. ZENG, 2002 Modeling epistasis of quantitative trait loci using cockerham's model. Genetics **160**: 1243–1261.

KAO, C. H., Z. B. ZENG, and R. D. TEASDALE, 1999 Multiple interval mapping for quantitative trait loci. Genetics **152**: 1203–1216.

KAPETANIOS, G. and M. WEEKS, 2003 Non-nested models and the likelihood ratio statistic: a comparison of simulation and bootstrap based tests. Technical report, University of London.

KASS, R. E. and A. E. RAFTERY, 1995 Bayes factors. Journal of the American Statistical Association **90**: 773–795.

KNOTT, S. A., J. ELSEN, and C. HALEY, 1996 Methods for multiple marker mapping of quantitative trait loci in half-sib. Theor Appl Genet **93**: 71–80.

KNOTT, S. A. and C. S. HALEY, 2000 Multitrait least squares for quantitative trait loci detection. Genetics **156**: 899–911.

LANDER, E. S. and D. BOTSTEIN, 1989 Mapping mendelian factors underlying quantitative traits using rflp linkage maps. Genetics **121**: 185–199.

LANGE, C. and J. C. WHITTAKER, 2001 Mapping quantitative trait loci using generalized estimating equations. Genetics **159**: 1325–1337.

LEBRETON, C. M., P. M. VISSCHER, C. S. HALEY, A. SEMIKHODSKII, and S. A. QUARRIET, 1998 A nonparametric bootstrap method for testing close linkage vs. pleiotropy of coincident quantitative trait loci. Genetics **150**: 931–943.

LIN, D., 2005 An efficient monte carlo approach to assessing statistical significance in genomic studies. Bioinformatics **21**: 781–787.

LIU, B.-H., 1998 *Statistical genomics : linkage, mapping, and QTL analysis*. CRC Press.

LIU, J., J. M. MERCER, L. F. STAM, G. C. GIBSON, Z. B. ZENG, and C. C. LAURIE, 1996 Genetic analysis of a morphological shape difference in the male genitalia of drosophila simulans and d. mauritiana. Genetics **142**: 1129–1145.

LIU, T. and R. WU, 2009 A bayesian algorithm for functional mapping of dynamic complex traits. Algorithms **2**: 667–691.

LYNCH, M. and B. WALSH, 1997 *Genetics and Analysis of Quantitative Traits*. Sinauer.

MA, C.-X., G. CASELLA, and R. WU, 2002 Functional mapping of quantitative trait loci underlying the character process: A theoretical framework. Genetics **161**: 1751–1762.

MACKAY, I. and W. POWELL, 2007 Methods for linkage disequilibrium mapping in crops. TRENDS in Plant Science **12**: 57–63.

MACKAY, T. F., 1996 The nature of quantittative genetic variation revisited: Lessons from drosophila bristles. BioEssays **18**: 113–121.

MACKAY, T. F., 2001 Quantitative trait loci in drosophila. Nat Rev Genet **2**: 11–20.

MAIA, J. M., 2007 *Joint Analysis of Multiple Gene Expression Traits to Map EXpression Quantitative Trait Loci*. Ph.D. thesis, North Carolina State University.

MANGIN, B., B. GOFFINET, and A. REBAI, 1994 Constructing confidence intervals for qtl location. Genetics **138**: 1301–1308.

MANGIN, B., P. THOQUET, and N. GRIMSLEY, 1998 Pleiotropic qtl analysis. Biometrics **54**: 88–99.

MCCARTHY, M. I., G. R. ABECASIS, L. R. CARDON, D. B. GOLDSTEIN, J. LITTLE, J. P. A. IOANNIDIS, and J. N. HIRSCHHORN, 2008 Genome-wide association studies for complex traits: consensus, uncertainty and challenges. Nat Rev Genet **9**: 356–369.

MCLACHLAN, G. J. and T. KRISHNAN, 1996 *The EM Algorithm and Extensions*. Wiley.

MENG, X.-L. and D. B. RUBIN, 1993 Maximum likelihood estimation via the ecm algorithm: A general framework. Biometrika **80**: 267–278.

METROPOLIS, N., A. W. ROSENBLUTH, M. N. ROSENBLUTH, and A. H. TELLER, 1953 Equation of state calculations by fast computing machines. The Journal of Chemical Physics **21**: 1087–1092.

MILLER, A., 2002 *Subset selection in regression*. Chapman & Hall/CRC.

NEALE, B., M. FERREIRA, S. MEDLAND, and D. POSTHUMA, editors, 2007 *Statistical Genetics: Gene Mapping Through Linkage and Association*. Taylor and Francis.

PATERSON, A. H., E. S. LANDER, J. D. HEWITT, S. PETERSON, S. E. LINCOLN, and S. D. TANKSLEY, 1988 Resolution of quantitative traits into mendelian factors by using a complete linkage map of restriction fragment length polymorphism. Nature **335**: 721–726.

RAI, S. N. and D. E. MATTHEWS, 1993 Improving the em algorithm. Biometrics **49**: 587–591.

RAO, C., 1948 Large sample testes of statistical hypothesis concerning several parameters with application to problems of estimation. Proceedings of the Cambridge Philosophical Society **44**: 50–57.

REBAI, A., B. GOFFINET, and B. MANGIN, 1994 Approximate thresholds of interval mapping tests for qtl detection. Genetics **138**: 235–240.

REBAI, A., B. GOFFINET, and B. MANGIN, 1995 Comparing power of different methods for qtl detection. Biometrics **51**: 87–99.

REDNER, R. A. and H. R. WALKER, 1984 Mixture densities, maximum likelihood and the em algorithm. SIAM Review **26**: 195–239.

REIFSNYDER, P. C., G. CHURCHILL, and E. H. LEITER, 2000 Maternal environment and genotype interact to establish diabesity in mice. Genome Res **10**: 1568–1578.

RENCHER, A. C., 2002 *Methods of Multivariate Analysis*. Wiley, second edition.

ROBERTS, S. B., C. J. MacLEAN, M. C. NEALE, L. J. EAVES, and K. S. KENDLER, 1999 Replication of linkage studies of complex traits: An examination of variation in location estimates. Am. J. Hum. Genet. **65**: 876–884.

SATAGOPAN, J. M. and B. S. YANDELL, 1996 Estimating the number of quantitative trait loci via bayesian model determination. In *Proceedings of the Joint Statistical Meetings, Chicago, IL, USA*, Special Contributed Paper Session on Genetic Analysis of Quantitative Traits and Complex Diseases, Biometric Section (available at http://www.stat.wisc.edu/ yandell/doc/1996/revjump/).

SATAGOPAN, J. M., B. S. YANDELL, M. A. NEWTON, and T. C. OSBORN, 1996 A bayesian approach to detect quantitative trait loci using markov chain monte carlo. Genetics **144**: 805–816.

SAX, K., 1923 Association of size differences with seed-coat pattern and pigmentation in phaseolus vulgaris. Genetics **8**: 552–560.

SCHWARZ, G., 1978 Estimating the dimension of a model. The Annals of Statistics **6**: 461–464.

SEN, S. and G. A. CHURCHILL, 2001 A statistical framework for quantitative trait mapping. Genetics **159**: 371–387.

SHAO, J. and D. TU, 1995 *The Jackknife and Bootstrap*. Springer.

SHRIMPTON, A. and ROBERTSON, 1988 The isolation of polygenic factors controlling bristle score in drosophila melanogaster. ii. distribution of third chromosome bristle effects within chromosome sections. Genetics **118**: 445–459.

SLATE, J., 2005 Quantitative trait locus mapping in natural populations: progress, caveats and future directions. Molecular Ecology **14**: 363–379.

SOLLER, M. and T. BRODY, 1976 On the power of experimental designs for the detection of linkage between marker loci and quantitative loci in crosses between inbred lines. Theor Appl Genet **47**: 35–39.

SORENSEN, D. and D. GIANOLA, 2002 *Likelihood, Bayesian, and MCMC methods in quantitative genetics*. Springer.

SUGIURA, N., 1978 Further analysts of the data by akaike' s information criterion and the finite corrections. Communications in Statistics - Theory and Methods **7**: 13–26.

TIBSHIRANI, R., 1996 Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society. Series B (Methodological) **58**: 267–288.

VAN DER BEEK, S., J. A. M. VAN ARENDONK, and A. F. GROEN, 1995 Power of two- and three-generation qtl mapping experiments in an outbred population containing full-sib or half-sib families. Theor Appl Genet **91**: 1115–1124.

VAN OOIJEN, J. W., 1992 Accuracy of mapping quantitative trait loci in autogamous species. Theoretical and Applied Genetics **84**: 803–811.

VISSCHER, P. M. and M. E. GODDARD, 2004 Prediction of the confidence interval of quantitative trait loci location. Behavior Genetics **34**: 477–482.

VISSCHER, P. M., R. THOMPSON, and C. S. HALEY, 1996 Confidence intervals in qtl mapping by bootstrapping. Genetics **143**: 1013–1020.

VUONG, Q. H., 1989 Likelihood ratio tests for model selection and non-nested hypotheses. Econometrica **57**: 307–333.

WAAGEPETERSEN, R. and D. SORENSEN, 2001 A tutorial on reversible jump mcmc with a view toward applications in qtl-mapping. International Statistical Review / Revue Internationale de Statistique **69**: 49–61.

WANG, H., Y.-M. ZHANG, X. LI, G. L. MASINDE, S. MOHAN, D. J. BAYLINK, and S. XU, 2005 Bayesian shrinkage estimation of quantitative trait loci parameters. Genetics **170**: 465–480.

WANG, S., C. J. BASTEN, and Z.-B. ZENG, 2007 Windows qtl cartographer 2.5. Department of Statistics, North Carolina State University, Raleigh, NC .

WEBER, K., R. EISMAN, L. MOREY, A. PATTY, J. SPARKS, M. TAUSEK, and Z.-B. ZENG, 1999 An analysis of polygenes affecting wing shape on chromosome 3 in drosophila melanogaster. Genetics **153**: 773786.

WELLER, J. I., G. WIGGANS, P. VANRADEN, and M. RON, 1996 Application of a canonical transformation to detection of quantitative trait loci with the aid of genetic markers in a multi-trait experiment. Theor Appl Genet **92**: 998–1002.

WU, R. and M. LIN, 2006 Functional mapping - how to map and study the genetic architecture of dynamic complex traits. Nat Rev Genet **7**: 229–237.

WU, R., C.-X. MA, and G. CASELLA, 2007 *Statistical genomics of quantitative traits: linkage, maps, and QTL*. Springer.

WU, W.-R., W.-M. LI, D.-Z. TANG, H.-R. LU, and A. J. WORLAND, 1999 Time-related mapping of quantitative trait loci underlying tiller number in rice. Genetics **151**: 297–303.

YANDELL, B. S., T. MEHTA, S. BANERJEE, D. SHRINER, R. VENKATARAMAN, J. Y. MOON, W. W. NEELY, H. WU, R. VON SMITH, and N. YI, 2007 R/qtlbim: Qtl with bayesian interval mapping in experimental crosses. Bioinformatics **23**: 641–643.

YANG, R., Q. TIAN, and S. XU, 2006 Mapping quantitative trait loci for longitudinal traits in line crosses. Genetics **173**: 2339–2356.

YANG, R. and S. XU, 2007 Bayesian shrinkage analysis of quantitative trait loci for dynamic traits. Genetics **176**: 1169–1185.

YAP, J. S., J. FAN, and R. WU, 2009 Nonparametric modeling of longitudinal covariance structure in functional mapping of quantitative trait loci. Biometrics **65**: 1068–1077.

YI, N., 2004 A unified markov chain monte carlo framework for mapping multiple quantitative trait loci. Genetics **167**: 967–975.

YI, N. and D. SHRINER, 2008 Advances in bayesian multiple quantitative trait loci mapping in experimental crosses. Heredity **100**: 240–252.

YI, N., D. SHRINER, S. BANERJEE, T. MEHTA, D. POMP, and B. S. YANDELL, 2007 An efficient bayesian model selection approach for interacting quantitative trait loci models with many effects. Genetics **176**: 1865–1877.

YI, N. and S. XU, 2008 Bayesian lasso for quantitative trait loci mapping. Genetics **179**: 1045–1055.

YI, N., S. XU, and D. B. ALLISON, 2003 Bayesian model choice and search strategies for mapping interacting quantitative trait loci. Genetics **165**: 867–883.

YI, N., B. S. YANDELL, G. A. CHURCHILL, D. B. ALLISON, E. J. EISEN, and D. POMP, 2005 Bayesian model selection for genome-wide epistatic quantitative trait loci analysis. Genetics **170**: 1333–1344.

ZELLNER, A., 1962 An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias. Journal of American Statistical Association **57**: 348–368.

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.

ZENG, Z. B., C. H. KAO, and C. J. BASTEN, 1999 Estimating the genetic architecture of quantitative traits. Genet Res **74**: 279–289.

ZENG, Z. B., J. LIU, L. F. STAM, C. H. KAO, J. M. MERCER, and C. C. LAURIE, 2000 Genetic architecture of a morphological shape difference between two drosophila species. Genetics **154**: 299–310.

ZENG, Z.-B., T. WANG, and W. ZOU, 2005 Modeling quantitative trait loci and interpretation of models. Genetics **169**: 1711–1725.

ZHAO, W., Y. Q. CHEN, G. CASELLA, J. M. CHEVERUD, and R. WU, 2005 A non-stationary model for functional mapping of complex traits. Bionformatics **21**: 2469–2477.

ZOU, F., J. P. FINE, J. HU, and D. Y. LIN, 2004 An efficient resampling method for assessing genome-wide statistical significance in mapping quantitative trait loci. Genetics **168**: 2307–2316.

ZOU, W. and Z.-B. ZENG, 2008 Statistical methods for mapping multiple qtl. Int J Plant Genomics **2008**: 286–561.

# Appendices

# A

# Introduction

## A.1  Joint and conditional probabilities of two QTL within a marker interval in backcross

Suppose two parental lines $P_1$ and $P_2$ with genotypes $m_L m_L q_1 q_1 q_2 q_2 m_R m_R$ and $M_L M_L Q_1 Q_1 Q_2 Q_2 M_R M_R$, respectively. Assume that the four loci are linked, i.e., they are relatively close to each other on a chromosome. Any offspring $F_1$ from the cross between $P_1$ and $P_2$ would have genotype $M_L m_L Q_1 q_1 Q_2 q_2 M_R m_R$. Any subject of a population obtained from the cross between $F_1$ and $P_2$ would have one of the sixteen genotypes shown in Table A.1. Their probabilities were computed assuming that the loci order in the chromosome is $M_L Q_1 Q_2 M_R$, and that the recombination frequencies between $M_L$ and $Q_1$, $Q_1$ and $Q_2$, $Q_2$ and $M_R$, and $M_L$ and $M_R$ are $r_1$, $r_2$, $r_3$, and $r$, respectively.

Table A.1: Joint probabilities of markers and QTL genotypes for any subject originated from the backcross of $F_1$ ($M_L m_L Q_1 q_1 Q_2 q_2 M_R m_R$) and $P_2$ ($M_L\ Q_1 Q_1 Q_2 Q_2 M_R M_R$), assuming the loci are linked.

| Genotype | Probability |
|---|---|
| $M_L M_L Q_1 Q_1 Q_2 Q_2 M_R M_R$ | $\frac{1}{2}[(1-r_1)(1-r_2)(1-r_2)-r_1 r_2(1-c_{12})-r_1 r_3(1-c_{13})-r_2 r_3(1-c_{23})+r_1 r_2 r_3(1-c)]$ |
| $M_L M_L Q_1 Q_1 Q_2 Q_2 M_R m_R$ | $\frac{1}{2}[(1-r_1)(1-r_2)r_3+r_1 r_3(1-c_{13})+r_2 r_3(1-c_{23})-r_1 r_2 r_3(1-c)]$ |
| $M_L M_L Q_1 Q_1 Q_2 q_2 M_R M_R$ | $\frac{1}{2}[(1-r_1)r_2 r_3-r_2 r_3(1-c_{23})+r_1 r_2 r_3(1-c)]$ |
| $M_L M_L Q_1 Q_1 Q_2 q_2 M_R m_R$ | $\frac{1}{2}[(1-r_1)r_2(1-r_3)+r_1 r_2(1-c_{12})+r_2 r_3(1-c_{23})-r_1 r_2 r_3(1-c)]$ |
| $M_L M_L Q_1 q_1 Q_2 Q_2 M_R M_R$ | $\frac{1}{2}[r_1 r_2(1-r_3)-r_1 r_2(1-c_{12})+r_1 r_2 r_3(1-c)]$ |
| $M_L M_L Q_1 q_1 Q_2 Q_2 M_R m_R$ | $\frac{1}{2}[c r_1 r_2 r_3]$ |
| $M_L M_L Q_1 q_1 Q_2 q_2 M_R M_R$ | $\frac{1}{2}[r_1(1-r_2)r_3-r_1 r_3(1-c_{13})+r_1 r_2 r_3(1-c)]$ |
| $M_L M_L Q_1 q_1 Q_2 q_2 M_R m_R$ | $\frac{1}{2}[r_1(1-r_2)(1-r_3)+r_1 r_2(1-c_{12})+r_1 r_3(1-c_{13})-r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 Q_1 Q_2 Q_2 M_R M_R$ | $\frac{1}{2}[r_1(1-r_2)(1-r_3)+r_1 r_2(1-c_{12})+r_1 r_3(1-c_{13})-r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 Q_1 Q_2 Q_2 M_R m_R$ | $\frac{1}{2}[r_1(1-r_2)r_3-r_1 r_2(1-c_{13})+r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 Q_1 Q_2 q_2 M_R M_R$ | $\frac{1}{2}[c r_1 r_2 r_3]$ |
| $M_L m_L Q_1 Q_1 Q_2 q_2 M_R m_R$ | $\frac{1}{2}[r_1 r_2(1-r_3)-r_1 r_2(1-c_{12})+r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 q_1 Q_2 Q_2 M_R M_R$ | $\frac{1}{2}[(1-r_1)r_2(1-r_3)+r_1 r_2(1-c_{12})+r_2 r_3(1-c_{23})-r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 q_1 Q_2 Q_2 M_R m_R$ | $\frac{1}{2}[(1-r_1)r_2 r_3-r_2 r_3(1-c_{23})+r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 q_1 Q_2 q_2 M_R M_R$ | $\frac{1}{2}[(1-r_1)(1-r_2)r_3+r_1 r_3(1-c_{13})+r_2 r_3(1-c_{23})-r_1 r_2 r_3(1-c)]$ |
| $M_L m_L Q_1 q_1 Q_2 q_2 M_R m_R$ | $\frac{1}{2}[(1-r_1)(1-r_2)(1-r_3)-r_1 r_2(1-c_{12})-r_1 r_3(1-c_{13})-r_2 r_3(1-c_{23})+r_1 r_2 r_3(1-c)]$ |

The order of markers and QTL is assumed to be $M_L Q_1 Q_2 M_R$, and the recombination frequencies between $M_L Q_1$, $Q_1 Q_2$, $Q_2 M_R$, and $M_L M_2$ are $r_1$, $r_2$, $r_3$, and $r$, respectively.

$c_{12} = \frac{\text{number of double recombinants } M_L \text{ and } Q_1, \text{ and } Q_1 \text{ and } Q_2}{r_1 r_2}$.

$c_{13} = \frac{\text{number of double recombinants } M_L \text{ and } Q_1, \text{ and } Q_2 \text{ and } M_R}{r_1 r_3}$.

$c_{23} = \frac{\text{number of double recombinants } Q_1 \text{ and } Q_2, \text{ and } Q_2 \text{ and } M_R}{r_2 r_3}$.

$c = \frac{\text{number of triple recombinants } M_L \text{and } Q_1, Q_1 \text{ and } Q_2, \text{ and } Q_2 \text{ and } M_R}{r_1 r_2 r_3}$.

Assuming the absence of double and triple cross-overs ($c_{12} = c_{13} = c_{23} = c = 0$), we derived the conditional probabilities of two QTL lying within a marker interval (Table A.2). The assumptions of no double and triple cross-overs greatly simplify the analytical derivation of probabilities.

Table A.2: Conditional probabilities of QTL genotypes $Q_1Q_1Q_2Q_2$, $Q_1Q_1Q_2q_2$, $Q_1q_1Q_2Q_2$, and $Q_1q_1Q_2q_2$, for any subject of a BC population, given the genotypes of markers $M_L$ (on the left) and $M_R$ (on the right) flanking the two QTL. The assumption of complete cross-over interference within a marker interval was made to obtain this table.

| Marker genotypes | $P(Q_1Q_1Q_2Q_2)$ | $P(Q_1Q_1Q_2q_2)$ | $P(Q_1q_1Q_2Q_2)$ | $P(Q_1q_1Q_2q_2)$ |
|---|---|---|---|---|
| $M_LM_LM_RM_R$ | 1 | 0 | 0 | 0 |
| $M_LM_LM_Rm_R$ | $\frac{r_3}{r}$ | $\frac{r_2}{r}$ | 0 | $\frac{r_1}{r}$ |
| $M_Lm_LM_RM_R$ | $\frac{r_1}{r}$ | 0 | $\frac{r_2}{r}$ | $\frac{r_3}{r}$ |
| $M_Lm_LM_Rm_R$ | 0 | 0 | 0 | 1 |

The order of markers and QTL is assumed to be $M_LQ_1Q_2M_R$, and the recombination frequencies between $M_LQ_1$, $Q_1Q_2$, $Q_2M_R$, and $M_LM_2$ are $r_1$, $r_2$, $r_3$, and $r$, respectively.

# B

# Multiple trait multiple interval mapping of quantitative trait loci from inbred line crosses

## B.1 First and second order derivatives of the logarithm likelihood function

In this section, we provide analytical formulae of the first and second order derivatives[*] of the logarithm of individual and overall likelihoods of MTMIM model. In what follows, for any matrix $\boldsymbol{A}$, its transpose is denoted by $\boldsymbol{A}'$, its inverse by $\boldsymbol{A}^{-1}$, its $u^{th}$ row by $\boldsymbol{A}_{[u,\cdot]}$, its $v^{th}$ column by $\boldsymbol{A}_{[\cdot,v]}$, and its element in row $u$ and column $v$ by $\boldsymbol{A}_{[u,v]}$.

In what follows, although all the definitions and notations have already been described throughout the chapters, we refer to them again in order to make this section self contented. Our MTMIM statistical model for QTL inference on BC population is a linear model, in which the value of trait $t$ ($t = 1, 2, \cdots, T$), $y_{ti}$, for

---

[*]We borrowed useful ideas from the works of DWYER and MACPHAIL (1948) and DWYER (1967). These papers provide many results regarding matrix derivatives as well their applications in multivariate analysis.

each $i^{th}$ subject ($i = 1, 2, \cdots, n$), is regressed on variables $x_{ir}$ ($r = 1, 2, \cdots, m$). These variables are defined according to the Cockerham genetic model (KAO and ZENG, 2002; ZENG *et al.*, 2005). For each subject $i$, $x_{ir}$ takes value $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether QTL $r$ has genotype $QQ$ or $Qq$, respectively. The coefficient of $x_{ir}$, $\beta_{tr}$, is called the main effect of $r^{th}$ QTL on trait $t$. The linear model also includes an intercept $\mu_t$, a subset $p$ of epistatic effects ($w_{trl}$) between all pairwise interactions between QTL ($r$ and $l \in \{1, 2, \cdots, m\}$), and a residue $e_{ti}$. The residues are assumed to be independent and identically distributed according to a normal distribution with mean zero and variance $\sigma_{e_t}^2$. The linear model is then:

$$y_{ti} = \mu_t + \sum_{r=1}^{m} \beta_{tr} x_{ir} + \sum_{r<l}^{p} w_{trl} x_{ir} x_{il} + e_{ti} \tag{B.1}$$

For each subject $i$, let $\boldsymbol{y}_i = (y_{1i}, y_{2i}, \cdots, y_{Ti})'$ be the $T$ by 1 vector of trait values, $\boldsymbol{X}_i$ be the $m + p$ by 1 incidence matrix, $\boldsymbol{e}_i$ be the $T$ by 1 vector of residuals, $\boldsymbol{\mu} = (\mu_1, \mu_2, \cdots, \mu_T)'$ be the $T$ by 1 vector of mean. For each $r$ and $l$, let $\boldsymbol{\beta}_r = (\beta_{1r}, \beta_{2r}, \cdots, \beta_{Tr})'$ and $\boldsymbol{w}_{rl} = (w_{1rl}, w_{2rl}, \cdots, w_{Trl})'$ be column vectors of main and epistatic effects, respectively. We collect all the effect parameters into a $T$ by $m+p$ matrix $\boldsymbol{\mathcal{B}} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \cdots, \boldsymbol{\beta}_m, \boldsymbol{w}_1, \boldsymbol{w}_2, \cdots, \boldsymbol{w}_p)$, and rewrite (B.1) in matrix form as:

$$\boldsymbol{y}_i = \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{X}_i + \boldsymbol{e}_i \tag{B.2}$$

where, $\boldsymbol{e}_i \sim MVN(\boldsymbol{0}, \boldsymbol{\Sigma}_e)$ with $\boldsymbol{\Sigma}_e$ being a positive definite symmetric matrix.

We collect all $s = m + p$ vectors of effect parameters ($m$ main and $p$ epistatic effect vectors), $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_e$ of model (B.2) into a vector $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}', vect(\boldsymbol{\Sigma}_e))'$, where $\boldsymbol{\theta}_b = \boldsymbol{\beta}_b'$ for $1 \leq b \leq m$ and $\boldsymbol{\theta}_b = \boldsymbol{w}_b'$ for $m < b \leq s$, and $vect(\boldsymbol{\Sigma}_e)$ is an operator that stacks the rows of $\boldsymbol{\Sigma}_e$ into a column vector one on the top of the other and then transposes it.

In the BC population there are two possible genotypes for each QTL, $QQ$ and $Qq$. Therefore, if there are $m$ QTL affecting a trait, there are $2^m$ possible genotypes

for any subject $i$. Genotypes of the form $G_j = Q_1Q_2 \cdots Q_m$, where $Q_r \in \{QQ, Qq\}$, $r = 1, 2, \cdots, m$ and $j = 1, 2, \cdots, 2^m$. We define an $s$ by $2^m$ matrix $\boldsymbol{Z}$ of coded genotypes according to the Cockerham genetic model (KAO and ZENG, 2002; ZENG *et al.*, 2005). In the matrix $\boldsymbol{Z}$ each row $b$ corresponds to a column of effect parameters in $\boldsymbol{\mathcal{B}}$ ($b = 1, 2, \cdots, s$) and each column $j$, $\boldsymbol{Z}_{[\cdot,j]}$, represents a coded genotype $G_j$. If $b \leq m$, $\boldsymbol{Z}_{[b,j]} = x_r$, otherwise $\boldsymbol{Z}_{[b,j]} = x_r * x_l$, where $x_u$ ($u = r$ or $u = l$) is either $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether the genotype of QTL $Q_u$ in $G_j$ is $QQ$ or $Qq$, respectively.

The natural logarithm of the individual ($\ell_i$) and overall likelihoods ($\ell$) are:

$$\ell_i(\boldsymbol{\theta}|\boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}) = \log_e\left(\sum_{j=1}^{2^m} p_{ij}(2\pi)^{-\frac{T}{2}}|\boldsymbol{\Sigma}_e|^{-\frac{1}{2}}e^{-\frac{1}{2}(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}Z}_{[\cdot,j]})'\boldsymbol{\Sigma}_e^{-1}(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}Z}_{[\cdot,j]})}\right)$$

$$\ell(\boldsymbol{\theta}|\boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}) = \sum_{i=1}^{n} \ell_i(\boldsymbol{\theta}|\boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda})$$

where, $p_{ij} = P(G_j|\boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\mathcal{R}}, \boldsymbol{\lambda})^{\dagger}$ is the conditional probability of QTL genotype $G_j = Q_1Q_2 \cdots Q_m$, where $Q_r \in \{QQ, Qq\}$, $r = 1, 2, \cdots, m$ and $j = 1, 2, \cdots, 2^m$.

For each subject $i$, let $\pi_{ij}$ be the posterior probability of QTL genotype $G_j$:

$$\pi_{ij} = \frac{p_{ij}\phi\left(\boldsymbol{y}_i|\boldsymbol{\mu} + \boldsymbol{\mathcal{B}Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e\right)}{\sum\limits_{j=1}^{2^m} p_{ij}\phi\left(\boldsymbol{y}_i|\boldsymbol{\mu} + \boldsymbol{\mathcal{B}Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e\right)}$$

where, $\phi(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ is the probability distribution function of a multivariate normal random variable $\boldsymbol{z}$ with mean $\boldsymbol{\mu}_0$ and variance-covariance $\boldsymbol{\Sigma}_0$.

*Auxiliary matrices*: we assume $b = 1, 2, \cdots, s$, $i = 1, 2, \cdots, n$ and $j = 1, 2, \cdots, 2^m$.

$\boldsymbol{J}$ is a $T$ by $T$ matrix with 1 at positions $\boldsymbol{J}_{[u,\ell]}$ and $\boldsymbol{J}_{[\ell,u]}$, and zero elsewhere.

$\boldsymbol{I}$ is an $T$ by $T$ identity matrix

---

†**Note on computations**: It is worth mentioning that for many combinations of $i$ and $j$, the probabilities $p_{ij}$ are zero or very close to zero. Therefore, one may take advantage of sparse matrix theory to save on computations when evaluating the expressions presented in Appendix B.

$$\boldsymbol{T}_{ij} = \boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]}$$

$$\boldsymbol{S}_{ij} = \tfrac{1}{2}\boldsymbol{T}_{ij}\boldsymbol{T}'_{ij}\boldsymbol{\Sigma}_e^{-1} - \tfrac{1}{2}\boldsymbol{I}$$

$$\frac{\partial \pi_{ij}}{\partial \boldsymbol{\theta}_b} = \pi_{ij}\boldsymbol{\Sigma}_e^{-1}\left(\boldsymbol{T}_{ij}\boldsymbol{Z}_{[b,j]} - \sum_{u=1}^{2^m} \pi_{iu}\boldsymbol{T}_{iu}\boldsymbol{Z}_{[b,u]}\right)$$

$$\frac{\partial \pi_{ij}}{\partial \boldsymbol{\mu}} = \pi_{ij}\boldsymbol{\Sigma}_e^{-1}\left(\boldsymbol{T}_{ij} - \sum_{u=1}^{2^m} \pi_{iu}\boldsymbol{T}_{iu}\right)$$

$$\frac{\partial \pi_{ij}}{\partial \boldsymbol{\Sigma}_e} = \pi_{ij}\boldsymbol{\Sigma}_e^{-1}\left(\boldsymbol{S}_{ij} - \sum_{u=1}^{2^m} \pi_{iu}\boldsymbol{S}_{iu}\right)$$

$$\frac{\partial \boldsymbol{T}_{ij}}{\partial \boldsymbol{\theta}'_b} = -\boldsymbol{Z}_{[b,j]}\boldsymbol{I}$$

$$\frac{\partial \boldsymbol{T}_{ij}}{\partial \boldsymbol{\mu}'} = -\boldsymbol{I}$$

$$\frac{\partial \phi\left(\boldsymbol{y}_i|\boldsymbol{\mu}+\boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]},\boldsymbol{\Sigma}_e\right)}{\partial \boldsymbol{\theta}_b} = \phi\left(\boldsymbol{y}_i|\boldsymbol{\mu}+\boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]},\boldsymbol{\Sigma}_e\right)\boldsymbol{\Sigma}_e^{-1}\boldsymbol{T}_{ij}\boldsymbol{Z}_{[b,j]}$$

$$\frac{\partial \phi\left(\boldsymbol{y}_i|\boldsymbol{\mu}+\boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]},\boldsymbol{\Sigma}_e\right)}{\partial \boldsymbol{\mu}} = \phi\left(\boldsymbol{y}_i|\boldsymbol{\mu}+\boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]},\boldsymbol{\Sigma}_e\right)\boldsymbol{\Sigma}_e^{-1}\boldsymbol{T}_{ij}$$

$$\frac{\partial \phi\left(\boldsymbol{y}_i|\boldsymbol{\mu}+\boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]},\boldsymbol{\Sigma}_e\right)}{\partial \boldsymbol{\Sigma}_e} = \phi\left(\boldsymbol{y}_i|\boldsymbol{\mu}+\boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]},\boldsymbol{\Sigma}_e\right)\boldsymbol{\Sigma}_e^{-1}\boldsymbol{S}_{ij}$$

*First order derivatives of the logarithm of the individual likelihood*:

In the following equations we short write $\ell_i(\boldsymbol{\theta}) = \ell_i(\boldsymbol{\theta}|\boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda})$, and we assume $b = 1, 2, \cdots, s$.

$$\frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_b} = \sum_{j=1}^{2^m} \pi_{ij}\boldsymbol{\Sigma}_e^{-1}\boldsymbol{T}_{ij}\boldsymbol{Z}_{[b,j]}$$

$$\frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}} = \sum_{j=1}^{2^m} \pi_{ij}\boldsymbol{\Sigma}_e^{-1}\boldsymbol{T}_{ij}$$

$$\frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_e} = \sum_{j=1}^{2^m} \pi_{ij}\boldsymbol{\Sigma}_e^{-1}\boldsymbol{S}_{ij}$$

*Second order derivatives of the logarithm of the overall likelihood*:

In the following equations we short write $\ell(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta}|\boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda})$, and we assume

$b = 1, 2, \cdots, s,\ k = 1, 2, \cdots, s,\ u = 1, 2, \cdots, T,\ \text{and}\ \ell = 1, 2, \cdots, T.$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_{k'} \partial \boldsymbol{\theta}_b} = \boldsymbol{\Sigma}_e^{-1} \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]} \boldsymbol{Z}_{[k,j]} \boldsymbol{T}_{ij} \boldsymbol{T}'_{ij} \right) \boldsymbol{\Sigma}_e^{-1}$$
$$- \boldsymbol{\Sigma}_e^{-1} \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]} \boldsymbol{T}_{ij} \sum_{c=1}^{2^m} \pi_{ic} \boldsymbol{Z}_{[k,c]} \boldsymbol{T}'_{ic} \right) \boldsymbol{\Sigma}_e^{-1}$$
$$- \boldsymbol{\Sigma}_e^{-1} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]} \boldsymbol{Z}_{[k,j]}$$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\mu}} = \boldsymbol{\Sigma}_e^{-1} \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{T}_{ij} \boldsymbol{T}'_{ij} \right) \boldsymbol{\Sigma}_e^{-1}$$
$$- \boldsymbol{\Sigma}_e^{-1} \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{T}_{ij} \sum_{c=1}^{2^m} \pi_{ic} \boldsymbol{T}'_{ic} \right) \boldsymbol{\Sigma}_e^{-1}$$
$$- n \boldsymbol{\Sigma}_e^{-1}$$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\theta}_b} = \boldsymbol{\Sigma}_e^{-1} \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]} \boldsymbol{T}_{ij} \boldsymbol{T}'_{ij} \right) \boldsymbol{\Sigma}_e^{-1}$$
$$- \boldsymbol{\Sigma}_e^{-1} \left( \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]} \boldsymbol{T}_{ij} \sum_{c=1}^{2^m} \pi_{ic} \boldsymbol{T}'_{ic} \right) \boldsymbol{\Sigma}_e^{-1}$$
$$- \boldsymbol{\Sigma}_e^{-1} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]}$$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}} \partial \boldsymbol{\theta}_b} = \sum_{i=1}^{n} \sum_{j=1}^{2^m} \frac{\partial \pi_{ij}}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}}} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij} \boldsymbol{Z}_{[b,j]}$$
$$- \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij} \boldsymbol{Z}_{[b,j]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij}$$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}} \partial \boldsymbol{\mu}} = \sum_{i=1}^{n} \sum_{j=1}^{2m} \frac{\partial \pi_{ij}}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}}} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij}$$
$$- \sum_{i=1}^{n} \sum_{j=1}^{2m} \pi_{ij} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij}$$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}} \partial \boldsymbol{\Sigma}_e} = \sum_{i=1}^{n} \sum_{j=1}^{2m} \frac{\partial \pi_{ij}}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}}} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{S}_{ij}$$
$$+ \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{2m} \pi_{ij} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1}$$
$$- \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{2m} \pi_{ij} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij} \boldsymbol{T}_{ij}' \boldsymbol{\Sigma}_e^{-1}$$
$$- \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{2m} \pi_{ij} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij} \boldsymbol{T}_{ij}' \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1}$$

*The matrices of first $(\frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}})$ and second $(\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}' \partial \boldsymbol{\theta}})$ order derivatives are:*

$$\frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \left( \begin{array}{ccccccc} \frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_1} & \frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_2} & \cdots & \frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_s} & | & \frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}} & | & \frac{\partial \ell_i(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_e} \end{array} \right)'$$

$$= \left( \begin{array}{ccccccc} \boldsymbol{C}_1 & \boldsymbol{C}_2 & \ldots & \boldsymbol{C}_s & | & \boldsymbol{V} & | & \boldsymbol{R} \end{array} \right)'$$

$$\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}' \partial \boldsymbol{\theta}} =$$

$$
\begin{pmatrix}
\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_1 \partial \boldsymbol{\theta}_1} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_s \partial \boldsymbol{\theta}_1} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\theta}_1} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{11}} \partial \boldsymbol{\theta}_1} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{TT}} \partial \boldsymbol{\theta}_1} \\
\vdots & \ddots & \vdots & \Big| & \vdots & \Big| & \vdots & \ddots & \vdots \\
\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_1 \partial \boldsymbol{\theta}_s} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_s \partial \boldsymbol{\theta}_s} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\theta}_s} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{11}} \partial \boldsymbol{\theta}_s} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{TT}} \partial \boldsymbol{\theta}_s} \\
\hline
\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_1 \partial \boldsymbol{\mu}} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_s \partial \boldsymbol{\mu}} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\mu}} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{11}} \partial \boldsymbol{\mu}} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{TT}} \partial \boldsymbol{\mu}} \\
\hline
\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_1 \partial \boldsymbol{\Sigma}_{e_{11}}} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_s \partial \boldsymbol{\Sigma}_{e_{11}}} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\Sigma}_{e_{11}}} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{11}} \partial \boldsymbol{\Sigma}_{e_{11}}} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{11}} \partial \boldsymbol{\Sigma}_{e_{TT}}} \\
\vdots & \ddots & \vdots & \Big| & \vdots & \Big| & \vdots & \ddots & \vdots \\
\frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_1 \partial \boldsymbol{\Sigma}_{e_{TT}}} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}'_s \partial \boldsymbol{\Sigma}_{e_{TT}}} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\Sigma}_{e_{TT}}} & \Big| & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{TT}} \partial \boldsymbol{\Sigma}_{e_{11}}} & \cdots & \frac{\partial^2 \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\Sigma}_{e_{TT}} \partial \boldsymbol{\Sigma}_{e_{TT}}}
\end{pmatrix}
$$

$$= \left( \begin{array}{ccccccccc}
\boldsymbol{A}_{1,1} & \dots & \boldsymbol{A}_{1,s} & & | & \boldsymbol{B}_1 & | & \boldsymbol{W}_{1,11} & \dots & \boldsymbol{W}_{1,TT} \\
\vdots & \ddots & \vdots & & | & \vdots & | & \vdots & \ddots & \vdots \\
\boldsymbol{A}'_{1,s} & \dots & \boldsymbol{A}_{s,s} & & | & \boldsymbol{B}_s & | & \boldsymbol{W}_{s,11} & \dots & \boldsymbol{W}_{s,TT} \\
\hline
\boldsymbol{B}'_1 & \dots & \boldsymbol{B}'_s & & | & \boldsymbol{K} & | & \boldsymbol{Y}_{11} & \dots & \boldsymbol{Y}_{TT} \\
\hline
\boldsymbol{W}'_{1,11} & \dots & \boldsymbol{W}'_{s,11} & & | & \boldsymbol{Y}'_{11} & | & \boldsymbol{Z}_{11,11} & \dots & \boldsymbol{Z}_{11,TT} \\
\vdots & \ddots & \vdots & & | & \vdots & | & \vdots & \ddots & \vdots \\
\boldsymbol{W}'_{1,TT} & \dots & \boldsymbol{W}'_{s,TT} & & | & \boldsymbol{Y}'_{TT} & | & \boldsymbol{Z}'_{11,TT} & \dots & \boldsymbol{Z}_{TT,TT}
\end{array} \right)$$

$$= \left( \begin{array}{ccccc}
\boldsymbol{A} & | & \boldsymbol{B} & | & \boldsymbol{W} \\
\hline
\boldsymbol{B}' & | & \boldsymbol{K} & | & \boldsymbol{Y} \\
\hline
\boldsymbol{W}' & | & \boldsymbol{Y}' & | & \boldsymbol{Z}
\end{array} \right)$$

## B.2 Score statistic

In this section, we derive the score statistic for testing any effect parameter in the MTMIM model (B.1). Assume model (B.1) with $m$ main and $p$ epistatic vectors, and collect all parameters into a column vector $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}', vect(\boldsymbol{\Sigma}_e))'$, where $s = m + p$. Assume that the hypotheses of interest are $\mathrm{H}_0 : \boldsymbol{\theta}_b = \boldsymbol{0}$ versus $\mathrm{H}_1 : \boldsymbol{\theta}_b \neq \boldsymbol{0}$

for some $b \in \{1, 2, \cdots, s\}$. Let $\boldsymbol{\eta} = (\boldsymbol{\theta_1}, \boldsymbol{\theta_2}, \cdots, \boldsymbol{\theta_{b-1}}, \boldsymbol{\theta_{b+1}}, \cdots, \boldsymbol{\theta_s}, \boldsymbol{\mu}, vect(\boldsymbol{\Sigma_e}))'$ be a column vector of nuisance parameters. Let $\tilde{\boldsymbol{\eta}}$ be the MLE of $\boldsymbol{\eta}$ under $H_0$. Then the score statistic to test $H_0$ versus $H_1$ can be defined as $S = \hat{\boldsymbol{U}}' \hat{\boldsymbol{V}}^{-1} \hat{\boldsymbol{U}}$ (Zou *et al.*, 2004; Cox and Hinkley, 1974), where $\boldsymbol{U} = \sum_{i=1}^{n} \hat{\boldsymbol{U}}_i$, $\hat{\boldsymbol{V}} = \sum_{i=1}^{n} \hat{\boldsymbol{U}}_i \hat{\boldsymbol{U}}_i'$, and $\hat{\boldsymbol{U}}_i$ is:

$$
\hat{\boldsymbol{U}}_i = \frac{\partial \ell_i (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\theta_b}} \bigg|_{(\boldsymbol{\theta_b}=0, \boldsymbol{\eta}=\tilde{\boldsymbol{\eta}})} -
$$

$$
\frac{\partial \ell (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\theta_b} \partial \boldsymbol{\eta}'} \bigg|_{(\boldsymbol{\theta_b}=0, \boldsymbol{\eta}=\tilde{\boldsymbol{\eta}})} \left( \frac{\partial \ell (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}'} \bigg|_{(\boldsymbol{\theta_b}=0, \boldsymbol{\eta}=\tilde{\boldsymbol{\eta}})} \right)^{-1} \frac{\partial \ell_i (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\eta}} \bigg|_{(\boldsymbol{\theta_b}=0, \boldsymbol{\eta}=\tilde{\boldsymbol{\eta}})}
$$

After removing redundant rows and columns of the matrices $\boldsymbol{W}$, $\boldsymbol{Y}$, $\boldsymbol{Z}$ and $\boldsymbol{R}$ generated due to the assumption that $\boldsymbol{\Sigma_e}$ is symmetric, and denoting the new matrices as $\tilde{\boldsymbol{W}}$, $\tilde{\boldsymbol{Y}}$, $\tilde{\boldsymbol{Z}}$ and $\tilde{\boldsymbol{R}}$, respectively, the components of $\boldsymbol{U}_i$ in terms of matrix notation are:

$$
\frac{\partial \ell_i (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\theta_b}} = \boldsymbol{C}_b
$$

$$
\frac{\partial \ell (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\theta_b} \partial \boldsymbol{\eta}'} = \left( \begin{array}{ccccccccccc} \boldsymbol{A}_{b,1} & \dots & \boldsymbol{A}_{b-1,b-1} & \boldsymbol{A}_{b+1,b+1} & \dots & \boldsymbol{A}_{b,s} & |\boldsymbol{B}_b| & \tilde{\boldsymbol{W}}_{b,11} & \dots & \tilde{\boldsymbol{W}}_{b,TT} \end{array} \right)
$$

$$
\frac{\partial \ell_i (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\eta}'} = \left( \begin{array}{ccccccccc} \boldsymbol{C}_1 & \dots & \boldsymbol{C}_{b-1} & \boldsymbol{C}_{b+1} & \dots & \boldsymbol{C}_s & | & \boldsymbol{V} & | & \tilde{\boldsymbol{R}} \end{array} \right)'
$$

$$
\frac{\partial \ell (\boldsymbol{\theta_b}, \boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}'} = \left( \begin{array}{ccc} \tilde{\boldsymbol{A}} & | \; \tilde{\boldsymbol{B}} \; | & \tilde{\boldsymbol{W}} \\ \hline \tilde{\boldsymbol{B}}' & | \; \boldsymbol{K} \; | & \tilde{\boldsymbol{Y}} \\ \hline \tilde{\boldsymbol{W}}' & | \; \tilde{\boldsymbol{Y}}' \; | & \tilde{\boldsymbol{Z}} \end{array} \right)
$$

where,

$$\tilde{A} = \begin{pmatrix} A_{1,1} & \ldots & A_{1,b-1} & A_{1,b+1} & \ldots & A_{1,s} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ A_{b-1,1} & \ldots & A_{b-1,b-1} & A_{b-1,b+1} & \ldots & A_{b-1,s} \\ A_{b+1,1} & \ldots & A_{b+1,b-1} & A_{b+1,b+1} & \ldots & A_{b+1,s} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ A_{s,1} & \ldots & A_{s,b-1} & A_{s,b+1} & \ldots & A_{s,s} \end{pmatrix}$$

$$\tilde{B} = \begin{pmatrix} B_{1,1} \\ \vdots \\ B_{b-1,1} \\ B_{b+1,1} \\ \vdots \\ B_{s,1} \end{pmatrix}$$

$$\tilde{W} = \begin{pmatrix} W_{1,11} & W_{1,12} & \ldots & W_{1,TT} \\ \vdots & \ldots & \ldots & \vdots \\ W_{b-1,11} & W_{b-1,12} & \ldots & W_{b-1,TT} \\ W_{b+1,11} & W_{b+1,12} & \ldots & W_{b+1,TT} \\ \vdots & \ldots & \ldots & \vdots \\ W_{s,11} & W_{s,12} & \ldots & W_{s,TT} \end{pmatrix}$$

# B.3 First and second order derivatives of the expected complete-data logarithm likelihood

For a current value of $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}, vect(\boldsymbol{\Sigma}_e))'$, denoted $\boldsymbol{\theta}^{(\nu)}$, the first and second order derivatives of the expected complete-data logarithm likelihood are shown bellow. We assume $b = 1, 2, \cdots, s$, $k = 1, 2, \cdots, s$, $u = 1, 2, \cdots, T$ and $\ell = 1, 2, \cdots, T$.

$$\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\theta}_b} = \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij} \boldsymbol{Z}_{[b,j]}$$

$$\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\mu}} = \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij}$$

$$\frac{\partial Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\Sigma}_e} = \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{S}_{ij}$$

$$\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\theta}_k' \partial \boldsymbol{\theta}_b} = -\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{Z}_{[b,j]} \boldsymbol{Z}_{[k,j]}$$

$$\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\mu}} = -n \boldsymbol{\Sigma}_e^{-1}$$

$$\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\mu}' \partial \boldsymbol{\theta}_b} = -\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{Z}_{[b,j]}$$

$$\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}} \partial \boldsymbol{\theta}_b} = -\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{Z}_{[b,j]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij}$$

$$\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}} \partial \boldsymbol{\mu}} = -\sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij}$$

$$\frac{\partial^2 Q_c(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\nu)})}{\partial \boldsymbol{\Sigma}_{e_{[u,\ell]}} \partial \boldsymbol{\Sigma}_e} = \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1}$$

$$-\frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij} \boldsymbol{T}_{ij}' \boldsymbol{\Sigma}_e^{-1}$$

$$-\frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{2^m} \pi_{ij}^{(\nu)} \boldsymbol{\Sigma}_e^{-1} \boldsymbol{T}_{ij} \boldsymbol{T}_{ij}' \boldsymbol{\Sigma}_e^{-1} \boldsymbol{J}_{[u,\ell]} \boldsymbol{\Sigma}_e^{-1}$$

## B.4 Extension of score statistic to intercross $F_2$

Our MTMIM statistical model for QTL inference in intercross $F_2$ population is a linear model, in which the value of trait $t$ ($t = 1, 2, \cdots, T$), $y_{ti}$, for each $i^{th}$ subject ($i = 1, 2, \cdots, n$), is regressed on the explanatory variables $x_{ir}$ and $z_{ir}$ ($r = 1, 2, \cdots, m$). These variables are defined according to the Cockerham genetic model. For each subject $i$, $x_{ir}$ takes value 1, 0 or $-1$, depending on whether QTL $r$ has genotype $QQ$, $Qq$ or $qq$, respectively, and $z_{ir}$ takes value $-\frac{1}{2}$, $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether QTL $r$ has genotype $QQ$, $Qq$ or $qq$, respectively (KAO and ZENG, 2002). The coefficient of $x_{ir}$, $a_{rt}^*$, is the additive effect of the $r^{th}$ QTL, and the coefficient of $z_{ir}$, $d_{rt}^*$, is the dominance effect of the $r^{th}$ QTL. The linear model also includes an intercept $\mu_t$, the additive by additive epistatic effects $w_{rkt}^{aa}$ between QTL $r$ and $k$ for a subset $p$ of all pairwise interactions, the additive by dominant epistatic effects $w_{rkt}^{ad}$ between QTL $r$ and $k$ for a subset $s$ of all pairwise interactions, the dominant by additive epistatic effects $w_{rkt}^{da}$ between QTL $k$ and $r$ for a subset $u$ of all pairwise interactions, the dominant by dominant epistatic effects $w_{rkt}^{dd}$ between QTL $r$ and $k$ for a subset $v$ of all pairwise interactions, and the residue $e_{ti}$, which is assumed to be independent and identically distributed according to a normal distribution with mean zero and

variance $\sigma_{e_t}^2$. The linear model is then:

$$
\begin{aligned}
y_{ti} = \mu_t \ &+ \sum_{r=1}^{m} \left( a_{tr}^* x_{ir} + d_{tr}^* z_{ir} \right) \\
&+ \sum_{r<k}^{p} w_{trk}^{aa} x_{ir} x_{ik} \\
&+ \sum_{r<k}^{o} w_{trk}^{ad} x_{ir} z_{ik} \\
&+ \sum_{r<k}^{u} w_{trk}^{da} z_{ir} x_{ik} \\
&+ \sum_{r<k}^{v} w_{trk}^{dd} z_{ir} z_{ik} \\
&+ e_{ti}
\end{aligned}
\tag{B.3}
$$

For each subject $i$, let $\boldsymbol{y}_i = (y_{1i}, y_{2i}, \cdots, y_{Ti})'$ be the $T$ by 1 vector of trait values, $\boldsymbol{X}_i$ be the $(2m + p + o + u + v)$ by 1 incidence matrix, $\boldsymbol{e}_i$ be the $T$ by 1 vector of residuals, $\boldsymbol{\mu} = (\mu_1, \mu_2, \cdots, \mu_T)'$ be the $T$ by 1 vector of mean. For each $r$ and $k$, let $\boldsymbol{a}_r^* = (a_{1r}^*, a_{2r}^*, \cdots, a_{Tr}^*)'$, $\boldsymbol{w}_{rk}^{aa} = (w_{1rk}^{aa}, w_{2rk}^{aa}, \cdots, w_{Trk}^{aa})'$, $\boldsymbol{w}_{rk}^{ad} = (w_{1rk}^{ad}, w_{2rk}^{ad}, \cdots, w_{Trk}^{ad})'$, $\boldsymbol{w}_{rk}^{da} = (w_{1rk}^{da}, w_{2rk}^{da}, \cdots, w_{Trk}^{da})'$, $\boldsymbol{w}_{rk}^{dd} = (w_{1rk}^{dd}, w_{2rk}^{dd}, \cdots, w_{Trk}^{dd})'$. We collect all the effect parameters into a $T$ by $s = 2m + p + o + u + v$ matrix $\boldsymbol{\mathcal{B}}$:

$$
\begin{aligned}
\boldsymbol{\mathcal{B}} = \ &( \boldsymbol{a}_1^*, \boldsymbol{a}_2^*, \cdots, \boldsymbol{a}_m^*, \\
&\boldsymbol{d}_1^*, \boldsymbol{d}_2^*, \cdots, \boldsymbol{d}_m^*, \\
&\boldsymbol{w}_1^{aa}, \boldsymbol{w}_2^{aa}, \cdots, \boldsymbol{w}_p^{aa}, \\
&\boldsymbol{w}_1^{ad}, \boldsymbol{w}_2^{ad}, \cdots, \boldsymbol{w}_o^{ad}, \\
&\boldsymbol{w}_1^{da}, \boldsymbol{w}_2^{da}, \cdots, \boldsymbol{w}_u^{da}, \\
&\boldsymbol{w}_1^{dd}, \boldsymbol{w}_2^{dd}, \cdots, \boldsymbol{w}_v^{dd} )
\end{aligned}
$$

Then, the statistical model in matrix notation, for subject $i$, would look like:

$$
\boldsymbol{y}_i = \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{X}_i + \boldsymbol{e}_i
\tag{B.4}
$$

where, $\boldsymbol{e}_i$ is a random vector of length $T$ assumed to be independent and identically

distributed according to a multivariate normal distribution with mean vector zero and positive definite symmetric variance-covariance matrix $\boldsymbol{\Sigma}_e$ ($MVN_T(\mathbf{0}, \boldsymbol{\Sigma}_e)$).

We collect all $s$ vectors of effect parameters, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}_e$ of model (B.4) into a column vector $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_s, \boldsymbol{\mu}, vect(\boldsymbol{\Sigma}_e))'$. In the $F_2$ population there are three possible genotypes for each QTL, $QQ$, $Qq$ and $qq$. Therefore, if there are $m$ QTL in the MTMIM model, there are $3^m$ possible genotypes for any subject $i$, genotypes of the form $G_j = Q_1 Q_2 \cdots Q_m$, where $Q_r \in \{QQ, Qq, qq\}$, $r = 1, 2, \cdots, m$ and $j = 1, 2, \cdots, 3^m$. Let $\boldsymbol{Z}$ be an $s \times 3^m$ matrix of coded genotypes, where each row $b$ corresponds to a column of effect parameters in $\boldsymbol{\mathcal{B}}$ ($b = 1, 2, \cdots, s$) and each column $j$, $\boldsymbol{Z}_{[\cdot,j]}$, represents a coded genotype $G_j$. Each element of $\boldsymbol{Z}$ is defined as follows:

$$
\boldsymbol{Z}_{[b,j]} = \begin{cases}
x_r & \text{if} \quad b \leq m \\
z_r & \text{if} \quad m < b \leq 2m \\
x_r x_k & \text{if} \quad 2m < b \leq 2m + p \\
x_r z_k & \text{if} \quad 2m + p < b \leq 2m + p + o \\
z_r x_k & \text{if} \quad 2m + p + o < b \leq 2m + p + o + u \\
z_r z_k & \text{if} \quad 2m + p + o + u < b \leq s
\end{cases}
$$

where, $x_r$ takes value 1, 0 or $-1$, depending on whether the genotype of $Q_r$ in $G_j$ is $QQ$, $Qq$ or $qq$, respectively, and $z_l$ ($l = r$ and $l = k$) takes value $-\frac{1}{2}$, $\frac{1}{2}$ or $-\frac{1}{2}$, depending on whether the genotype of $Q_l$ in $G_j$ is $QQ$, $Qq$ or $qq$, respectively.

The natural logarithm of the individual likelihood ($\ell_i$) is:

$$
\ell_i\left(\boldsymbol{\theta} \mid \boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right) = \log_e\left(\sum_{j=1}^{3^m} p_{ij} (2\pi)^{-\frac{T}{2}} |\boldsymbol{\Sigma}_e|^{-\frac{1}{2}} e^{-\frac{1}{2}\left(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]}\right)'\boldsymbol{\Sigma}_e^{-1}\left(\boldsymbol{y}_i - \boldsymbol{\mu} - \boldsymbol{\mathcal{B}}\boldsymbol{Z}_{[\cdot,j]}\right)}\right)
$$

(B.5)

where, $p_{ij} = p(G_j | \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\mathcal{R}}, \boldsymbol{\lambda})$ is the conditional probability of QTL genotype $G_j = Q_1 Q_2 \cdots Q_m$, where $Q_r \in \{QQ, Qq, qq\}$, $r = 1, 2, \cdots, m$ and $j = 1, 2, \cdots, 3^m$. The

overall likelihood ($\ell$) is:

$$\ell\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{M}, \boldsymbol{\lambda}\right) = \sum_{i=1}^{n} \ell_i \left(\boldsymbol{\theta} \mid \boldsymbol{y}_i, \boldsymbol{M}_{[i,\cdot]}, \boldsymbol{\lambda}\right) \tag{B.6}$$

For each subject $i$, let $\pi_{ij}$ be the posterior probability of the QTL genotype $G_j$:

$$\pi_{ij} = \frac{p_{ij}\phi\left(\boldsymbol{y}_i \mid \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e\right)}{\sum\limits_{j=1}^{3^m} p_{ij}\phi\left(\boldsymbol{y}_i \mid \boldsymbol{\mu} + \boldsymbol{\mathcal{B}} \boldsymbol{Z}_{[\cdot,j]}, \boldsymbol{\Sigma}_e\right)}$$

where, $\phi(\boldsymbol{z} \mid \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ is the probability density distribution of a multivariate normal random variable $\boldsymbol{z}$ with mean $\boldsymbol{\mu}_0$ and variance-covariance $\boldsymbol{\Sigma}_0$.

The extension of score statistic to the $F_2$ population is straightforward, in fact, the auxiliary matrices, expressions of first and second order derivatives of the logarithm likelihood functions (B.5) and (B.6) with respect to the parameters in model (B.3) can be obtained straight from the general expressions derived in section B.1. The extension consists basically of using $\boldsymbol{Z}$ and $\boldsymbol{\mathcal{B}}$ matrices of the $F_2$ and substituting $2^m$ by $3^m$ in the summations where needed.

# C

# Prediction of length of confidence interval for position of QTL in multiple trait analysis

## C.1 Prediction of length of confidence interval for position of QTL

We assume that measurements are taken in $T$ traits for each subject $i$ of a BC population in which a QTL may assume genotype $QQ$ and $Qq$, with additive effect $a_t$ and dominance effect $d_t$ on trait $t$ ($t = 1, 2, \cdots, T$), respectively ((FALCONER and MACKAY, 1996)). Each QTL genotype has expected frequency of $\frac{1}{2}$. Similarly, any genetic marker has either genotype $MM$ or $Mm$, each with expected frequency of $\frac{1}{2}$. We define $\boldsymbol{a} = (a_1, a_2, \cdots, a_T)'$ and $\boldsymbol{d} = (d_1, d_2, \cdots, d_T)'$.

We define a random variable $X_t(G)$ that takes value either $a_t$ or $d_t$, according to whether $G = QQ$ or $G = Qq$, respectively, with probability $\frac{1}{2}$ for both outcomes. Then, expectation of $X_t(G)$ is $\frac{1}{2}(a_t + d_t)$, and expectation of $X_t^2(G)$ is $\frac{1}{2}(a_t^2 + d_t^2)$. We define another variable $X_{t'}(G)$ that also takes values either $a_{t'}$ or $d_{t'}$, according to whether $G = QQ$ or $G = Qq$, respectively, with probability $\frac{1}{2}$ for both outcomes.

With these two variables we obtain $\mathrm{E}[X_t(G)X_{t'}(G)] = \frac{1}{2}(a_t a_{t'} + d_t d_{t'})$, variance of $X_t(G)$, and covariance of $X_t(G)$ and $X_{t'}(G)$ as:

$$\begin{aligned}
\mathrm{Var}[X_t(G)] =& \mathrm{E}[X_t^2(G)] - \mathrm{E}^2[X_t(G)] \\
=& \frac{1}{4}(a_t - d_t)^2 \\
\mathrm{Cov}[X_t(G)X_{t'}(G)] =& \mathrm{E}[X_t(G)X_{t'}(G)] - \mathrm{E}[X_t(G)]\mathrm{E}[X_{t'}(G)] \\
=& \frac{1}{4}(a_t - d_t)(a_{t'} - d_{t'})
\end{aligned}$$

For any pair of genetic marker and QTL, and assuming a large population we are able to observe the four haplotypes $MQ$, $Mq$, $mQ$ and $mq$. Assuming that the genetic marker locus and the QTL are linked in *cis* phase (i.e., in the $F_1$ generation $M$ and $Q$ resides in the same chromosome) and have recombination frequency $r$, the four haplotypes have expected frequencies $\frac{1-r}{2}$, $\frac{r}{2}$, $\frac{r}{2}$, and $\frac{1-r}{2}$, respectively. For each trait $t$ and a given genetic marker we define the phenotypic mean of the four haplotypes as $\bar{y}_{t_{MQ}}$, $\bar{y}_{t_{Mq}}$, $\bar{y}_{t_{mQ}}$, and $\bar{y}_{t_{mq}}$, respectively. We collect the means of all traits for each haplotype into four vectors $\bar{\boldsymbol{Y}}_{MQ} = (\bar{y}_{1_{MQ}}, \bar{y}_{2_{MQ}}, \cdots, \bar{y}_{T_{MQ}})'$, $\bar{\boldsymbol{Y}}_{Mq} = (\bar{y}_{1_{Mq}}, \bar{y}_{2_{Mq}}, \cdots, \bar{y}_{T_{Mq}})'$, $\bar{\boldsymbol{Y}}_{mQ} = (\bar{y}_{1_{mQ}}, \bar{y}_{2_{mQ}}, \cdots, \bar{y}_{T_{mQ}})'$, and $\bar{\boldsymbol{Y}}_{mq} = (\bar{y}_{1_{mq}}, \bar{y}_{2_{mq}}, \cdots, \bar{y}_{T_{mq}})'$. We also define the phenotypic means of the two genotypes $MM$ and $Mm$ of the genetic marker as $\bar{y}_{t_M}$ and $\bar{y}_{t_m}$, respectively. We collect the means of genotypes $MM$ and $Mm$ for all traits into the vectors $\bar{\boldsymbol{Y}}_M = (\bar{y}_{1_M}, \bar{y}_{2_M}, \cdots, \bar{y}_{T_M})'$ and $\bar{\boldsymbol{Y}}_m = (\bar{y}_{1_m}, \bar{y}_{2_m}, \cdots, \bar{y}_{T_m})'$, respectively. Likewise, we collect the phenotypic means of genotypes $QQ$ and $Qq$ into the vectors $\bar{\boldsymbol{Y}}_Q = (\bar{y}_{1_Q}, \bar{y}_{2_Q}, \cdots, \bar{y}_{T_Q})'$ and $\bar{\boldsymbol{Y}}_q = (\bar{y}_{1_q}, \bar{y}_{2_q}, \cdots, \bar{y}_{T_q})'$, respectively. The phenotypic means of marker and QTL can be written is terms of haplotype means as

follows:

$$\bar{Y}_M = (1-r)\bar{Y}_{MQ} + r\bar{Y}_{Mq}$$

$$\bar{Y}_m = (1-r)\bar{Y}_{mq} + r\bar{Y}_{mQ}$$

$$\bar{Y}_Q = (1-r)\bar{Y}_{MQ} + r\bar{Y}_{mQ}$$

$$\bar{Y}_q = (1-r)\bar{Y}_{mq} + r\bar{Y}_{Mq}$$

Assuming that $\mathbf{\Sigma}_{(\bar{Y}_M - \bar{Y}_m)} = \mathbf{\Sigma}_{(\bar{Y}_Q - \bar{Y}_q)} \approx \frac{4}{n}\mathbf{\Sigma}_e$, where $\mathbf{\Sigma}_e$ is the residual variance-covariance matrix, we define a $D(l)$ statistic as:

$$
\begin{aligned}
D(l) =& (\bar{Y}_M - \bar{Y}_m)'\mathbf{\Sigma}^{-1}_{(\bar{Y}_M - \bar{Y}_m)}(\bar{Y}_M - \bar{Y}_m) - (\bar{Y}_Q - \bar{Y}_q)'\mathbf{\Sigma}^{-1}_{(\bar{Y}_Q - \bar{Y}_q)}(\bar{Y}_Q - \bar{Y}_q) \\
=& \frac{n}{4}[(\bar{Y}_M - \bar{Y}_m)'\mathbf{\Sigma}^{-1}_e(\bar{Y}_M - \bar{Y}_m) - (\bar{Y}_Q - \bar{Y}_q)'\mathbf{\Sigma}^{-1}_e(\bar{Y}_Q - \bar{Y}_q)]
\end{aligned}
$$

The statistic $D(l)$ can be expressed in terms of means of four haplotypes in the population as shown bellow. Let

$$
\begin{aligned}
(\bar{Y}_M - \bar{Y}_m)'\mathbf{\Sigma}^{-1}_e(\bar{Y}_M - \bar{Y}_m) =& [(1-r)\bar{Y}_{MQ} + r\bar{Y}_{Mq} - (1-r)\bar{Y}_{mq} - r\bar{Y}_{mQ}]'\mathbf{\Sigma}^{-1}_e \\
& [(1-r)\bar{Y}_{MQ} + r\bar{Y}_{Mq} - (1-r)\bar{Y}_{mq} - r\bar{Y}_{mQ}] \\
=& [(1-r)(\bar{Y}_{MQ} - \bar{Y}_{mq}) + r(\bar{Y}_{Mq} - \bar{Y}_{mQ})]'\mathbf{\Sigma}^{-1}_e \\
& [(1-r)(\bar{Y}_{MQ} - \bar{Y}_{mq}) + r(\bar{Y}_{Mq} - \bar{Y}_{mQ})]
\end{aligned}
$$

and

$$
\begin{aligned}
(\bar{Y}_Q - \bar{Y}_q)'\mathbf{\Sigma}^{-1}_e(\bar{Y}_Q - \bar{Y}_q) =& [(1-r)\bar{Y}_{MQ} + r\bar{Y}_{mQ} - (1-r)(\bar{Y}_{mq} - r\bar{Y}_{Mq})]'\mathbf{\Sigma}^{-1}_e \\
& [(1-r)\bar{Y}_{MQ} + r\bar{Y}_{mQ} - (1-r)(\bar{Y}_{mq} - r\bar{Y}_{Mq})] \\
=& [(1-r)(\bar{Y}_{MQ} - \bar{Y}_{mq}) - r(\bar{Y}_{Mq} - \bar{Y}_{mQ})]'\mathbf{\Sigma}^{-1}_e \\
& [(1-r)(\bar{Y}_{MQ} - \bar{Y}_{mq}) - r(\bar{Y}_{Mq} - \bar{Y}_{mQ})]
\end{aligned}
$$

then,

$$
\begin{aligned}
D(l) =& \frac{n}{4}\{[(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq}) + r(\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})]'\boldsymbol{\Sigma}_e^{-1} \\
& [(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq}) + r(\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})] \\
& - [(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq}) - r(\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})]'\boldsymbol{\Sigma}_e^{-1} \\
& [(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq}) - r(\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})]\} \\
=& \frac{n}{4}[4r(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq})'\boldsymbol{\Sigma}_e^{-1}(\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})] \\
=& nr(1-r)(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq})'\boldsymbol{\Sigma}_e^{-1}(\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})
\end{aligned}
$$

In order to obtain the expected value of $D(l)$, we derive some auxiliary equations (equations (C.1), (C.2) (C.3), (C.4) and (C.5)) as follows:

$$
\begin{aligned}
\mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m) =& \mathrm{E}(\bar{\boldsymbol{Y}}_M) - \mathrm{E}(\bar{\boldsymbol{Y}}_m) \\
=& (1-r)\boldsymbol{a} + r\boldsymbol{d} - r\boldsymbol{a} - (1-r)\boldsymbol{d} \\
=& (1-2r)\boldsymbol{a} - (1-2r)\boldsymbol{d} \\
=& (1-2r)(\boldsymbol{a} - \boldsymbol{d}) \tag{C.1}
\end{aligned}
$$

By setting $r = 0$ in equation (C.1), we obtain $\mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q) = \boldsymbol{a} - \boldsymbol{d}$. Under the assumption of additivity of allelic effects ($\boldsymbol{d} = \boldsymbol{0}$), $\mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)$ and $\mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)$ are:

$$
\mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m) = (1-2r)\boldsymbol{a} \tag{C.2}
$$

$$
\mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q) = \boldsymbol{a} \tag{C.3}
$$

The expectation of the cross-products are:

$$\mathrm{E}[(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)'] \approx \frac{4}{n}\boldsymbol{\Sigma}_e + \mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)\mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)'$$

$$= \frac{4}{n}\boldsymbol{\Sigma}_e + (1-2r)^2\boldsymbol{a}\boldsymbol{a}' \tag{C.4}$$

$$\mathrm{E}[(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)'] \approx \frac{4}{n}\boldsymbol{\Sigma}_e + \mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)\mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)'$$

$$= \frac{4}{n}\boldsymbol{\Sigma}_e + \boldsymbol{a}\boldsymbol{a}' \tag{C.5}$$

Using equations (C.4) and (C.5), and applying some properties of trace of matrices we obtain the desired expected value of $D(l)$ as follows:

$$\begin{aligned}
\mathrm{E}[D(l)] =& \frac{n}{4}\{\mathrm{E}[tr[(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)\boldsymbol{\Sigma}_e^{-1}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)']] - \mathrm{E}[tr[(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)\boldsymbol{\Sigma}_e^{-1}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)']]\} \\
=& \frac{n}{4}\{tr[\boldsymbol{\Sigma}_e^{-1}\mathrm{E}[(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)']] - tr[\boldsymbol{\Sigma}_e^{-1}\mathrm{E}[(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)']]\} \\
=& \frac{n}{4}\{tr[\boldsymbol{\Sigma}_e^{-1}[\frac{4}{n}\boldsymbol{\Sigma}_e + \mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)\mathrm{E}(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)']] \\
& - tr[\boldsymbol{\Sigma}_e^{-1}[\frac{4}{n}\boldsymbol{\Sigma}_e + \mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)\mathrm{E}(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)']]\} \\
=& \frac{n}{4}\{tr[\boldsymbol{\Sigma}_e^{-1}(\frac{4}{n}\boldsymbol{\Sigma}_e + (1-2r)^2\boldsymbol{a}\boldsymbol{a}')] - tr[\boldsymbol{\Sigma}_e^{-1}(\frac{4}{n}\boldsymbol{\Sigma}_e + \boldsymbol{a}\boldsymbol{a}')]\} \\
=& \frac{n}{4}\{\frac{4}{n}tr[\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}_e] + (1-2r)^2 tr[\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}\boldsymbol{a}'] - \frac{4}{n}tr[\boldsymbol{\Sigma}_e^{-1}\boldsymbol{\Sigma}_e] + tr[\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}\boldsymbol{a}']\} \\
=& \frac{n}{4}[-4r(1-r)tr(\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}\boldsymbol{a}')] \\
=& -nr(1-r)\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a} \tag{C.6}
\end{aligned}$$

The expectation of $D^\star(l) = \frac{D(l)}{nr(1-r)}$ is easily obtained from equation (C.6):

$$\mathrm{E}[D^\star(l)] = -\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}\boldsymbol{a}$$

In what follows, we derive auxiliary equations (equations (C.7), (C.8), (C.9),

(C.10), (C.11), (C.13) and (C.12)) to compute the variance of $D^\star(l)$:

$$
\begin{aligned}
\mathrm{E}(\bar{y}_{t_{MQ}} - \bar{y}_{t_{mq}}) &= \mathrm{E}(\bar{y}_{t_{MQ}}) - \mathrm{E}(\bar{y}_{t_{mq}}) \\
&= a_t - (-d_t) \\
&= a_t
\end{aligned}
\tag{C.7}
$$

$$
\begin{aligned}
\mathrm{E}(\bar{y}_{t_{Mq}} - \bar{y}_{t_{mQ}}) &= \mathrm{E}(\bar{y}_{t_{Mq}}) - \mathrm{E}(\bar{y}_{t_{mQ}}) \\
&= -d_t - a_t \\
&= -a_t
\end{aligned}
\tag{C.8}
$$

$$
\mathrm{Cov}(\bar{y}_{t_{MQ}} - \bar{y}_{t_{mq}}, \bar{y}_{t'_{MQ}} - \bar{y}_{t'_{mq}}) = \frac{4\sigma_{e_{tt'}}}{n(1-r)}
\tag{C.9}
$$

$$
\mathrm{Cov}(\bar{y}_{t_{Mq}} - \bar{y}_{t_{mQ}}, \bar{y}_{t'_{Mq}} - \bar{y}_{t'_{mQ}}) = \frac{4\sigma_{e_{tt'}}}{nr}
\tag{C.10}
$$

In equations (C.7) and (C.8), we assumed $d_t = 0$ (additive allelic effect). In equations (C.9) and (C.10), $\sigma_{e_{tt'}}$ is the element in row $t$ and column $t'$ of $\boldsymbol{\Sigma}_e$. The last two equations above were obtained under the assumption of $\boldsymbol{\Sigma}_{(\bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m)} = \boldsymbol{\Sigma}_{(\bar{\boldsymbol{Y}}_Q - \bar{\boldsymbol{Y}}_q)} \approx \frac{4}{n}\boldsymbol{\Sigma}_e$.

For any variables $X, Y, W$ and $Z$, if $X$ and $W$ are both independents of $Y$ and $Z$, and $\rho_{XW}$ and $\rho_{YZ}$ are small, then $\mathrm{Cov}(XY, WZ) \approx \mathrm{Cov}(XW)\mathrm{E}(Y)\mathrm{E}(Z) +$

$E(X)(W)\mathrm{Cov}(Y,Z)$, as shown bellow:

$$
\begin{aligned}
\mathrm{Cov}(XY,WZ) =& \mathrm{E}(XYWZ) - \mathrm{E}(XY)\mathrm{E}(WZ) \\
=& \mathrm{E}(XYWZ) - \mathrm{E}(X)\mathrm{E}(Y)\mathrm{E}(W)\mathrm{E}(Z) \\
=& \mathrm{E}(XW)\mathrm{E}(YZ) - \mathrm{E}(X)\mathrm{E}(Y)\mathrm{E}(W)\mathrm{E}(Z) \\
=& (\mathrm{Cov}(XW) + \mathrm{E}(X)\mathrm{E}(W))\mathrm{E}(YZ) - \mathrm{E}(X)\mathrm{E}(Y)\mathrm{E}(W)\mathrm{E}(Z) \\
=& \mathrm{Cov}(XW)\mathrm{E}(YZ) + \mathrm{E}(X)(W)\mathrm{Cov}(Y,Z) \\
=& \mathrm{Cov}(XW)(\mathrm{Cov}(YZ) + \mathrm{E}(Y)\mathrm{E}(Z)) + \mathrm{E}(X)(W)\mathrm{Cov}(Y,Z) \\
=& \mathrm{Cov}(XW)\mathrm{E}(Y)\mathrm{E}(Z) + \mathrm{E}(X)(W)\mathrm{Cov}(Y,Z) + \mathrm{Cov}(X,W)\mathrm{Cov}(Y,Z) \\
=& \mathrm{Cov}(XW)\mathrm{E}(Y)\mathrm{E}(Z) + \mathrm{E}(X)(W)\mathrm{Cov}(Y,Z) + \rho_{XW}\rho_{YZ}\sigma_X\sigma_W\sigma_Y\sigma_Z \\
\approx& \mathrm{Cov}(XW)\mathrm{E}(Y)\mathrm{E}(Z) + \mathrm{E}(X)(W)\mathrm{Cov}(Y,Z) \qquad\qquad \text{(C.11)}
\end{aligned}
$$

Applying the results of equations (C.7), (C.8), (C.9), (C.10) and (C.11), where in the latter we let $X = (\bar{y}_{j_{MQ}} - \bar{y}_{j_{mq}})$, $Y = (\bar{y}_{k_{Mq}} - \bar{y}_{k_{mQ}})$, $W = (\bar{y}_{l_{MQ}} - \bar{y}_{l_{mq}})$, and $Z = (\bar{y}_{u_{Mq}} - \bar{y}_{u_{mQ}})$, we obtain:

$$
\mathrm{Cov}[(\bar{y}_{j_{MQ}} - \bar{y}_{j_{mq}})(\bar{y}_{k_{Mq}} - \bar{y}_{k_{mQ}}), (\bar{y}_{l_{MQ}} - \bar{y}_{l_{mq}})(\bar{y}_{u_{Mq}} - \bar{y}_{u_{mQ}})]
$$
$$
\approx \frac{4}{n}\left(\frac{a_j a_l \sigma_{e_{ku}}}{r} + \frac{a_k a_u \sigma_{e_{jl}}}{1-r}\right)
$$
$$
\text{(C.12)}
$$

where, $\sigma_{e_{jl}}$ is the element in row $j$ and column $l$ of the variance covariance-matrix $\boldsymbol{\Sigma}_e$ (likewise for $\sigma_{e_{ku}}$).

Let $x = (x_1, x_2, \cdots, x_T)'$, $y = (y_1, y_2, \cdots, y_T)'$ and $A$ be a square $T$ by $T$ matrix,

then the variance of the quadratic form $x'Ay$ is:

$$
\begin{aligned}
\mathrm{Var}(x'Ay) =& \mathrm{Var}[\sum_{t=1}^{T}\sum_{t'=1}^{T} a_{tt'} x_t y_{t'}] \\
=& \sum_{t=1}^{T}\sum_{t'=1}^{T} \mathrm{Var}[a_{tt'} x_t y_{t'}] \\
& + \sum_{t=1}^{T}\sum_{t'=1}^{T} \mathrm{Cov}[a_{tt'} x_t y_{t'}, a_{tt'} x_t y_{t'}]
\end{aligned}
\tag{C.13}
$$

Now, we let $x = (\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq})$, $y = (\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})$, $a_{tt'} = \sigma_{e_{tt'}}$ and $A = \boldsymbol{\Sigma}_e^{-1}$. Then from equations (C.12) and (C.13), we obtain the variance of $D^{\star}(l)$ as:

$$
\begin{aligned}
\mathrm{Var}[D^{\star}(l)] =& \mathrm{Var}[(\bar{\boldsymbol{Y}}_{MQ} - \bar{\boldsymbol{Y}}_{mq})' \boldsymbol{\Sigma}_e^{-1} (\bar{\boldsymbol{Y}}_{Mq} - \bar{\boldsymbol{Y}}_{mQ})] \\
=& \mathrm{Var}[\sum_{t=1}^{T}\sum_{t'=1}^{T} \sigma_{e_{tt'}} (\bar{y}_{t_{Mq}} - \bar{y}_{t_{mQ}})(\bar{y}_{t'_{Mq}} - \bar{Y}_{t'_{mQ}})] \\
=& \sum_{t=1}^{T}\sum_{t'=1}^{T} \mathrm{Var}[\sigma_{e_{tt'}} (\bar{y}_{t_{Mq}} - \bar{y}_{t_{mQ}})(\bar{y}_{t'_{Mq}} - \bar{y}_{t'_{mQ}})] \\
& + \sum_{t=1}^{T}\sum_{t'=1}^{T} \mathrm{Cov}[\sigma_{e_{tt'}} (\bar{y}_{t_{Mq}} - \bar{y}_{t_{mQ}})(\bar{y}_{t'_{Mq}} - \bar{y}_{t'_{mQ}}), \sigma_{e_{tt'}} (\bar{y}_{t_{Mq}} - \bar{y}_{t_{mQ}})(\bar{y}_{t'_{Mq}} - \bar{y}_{t'_{mQ}})] \\
=& \frac{4}{nr(1-r)} \boldsymbol{a}' \boldsymbol{\Sigma}_e^{-1} \boldsymbol{\Sigma}_e \boldsymbol{\Sigma}_e^{-1} \boldsymbol{a} \\
=& \frac{4}{nr(1-r)} \boldsymbol{a}' \boldsymbol{\Sigma}_e^{-1} \boldsymbol{a}
\end{aligned}
\tag{C.14}
$$

In terms of Haldane's distance $l$, $r = \frac{1}{2}(1 - e^{-2l})$, which leads to $r(1-r) = \frac{1}{4}(1 - e^{-4l})$. Thus the variance of $D^{\star}(l)$ in (C.14) can be rewritten as:

$$
\mathrm{Var}[D^{\star}(l)] = \frac{16}{n(1 - e^{-4l})} \boldsymbol{a}' \boldsymbol{\Sigma}_e^{-1} \boldsymbol{a}
$$

## C.2    Prediction of shape of LRT around the position of QTL

The phenotypic value of trait $t$ measured on subject $i$ ($i = 1, 2, \cdots, n$) may be written in the following linear model:

$$y_{ti} = \beta_t X_i + e_{ti} \tag{C.15}$$

where, $X_i$ takes value of $\frac{1}{2}$ or $-\frac{1}{2}$, according to whether subject $i$ has genotype $QQ$ or $Qq$, respectively, and $e_{ti}$ is the residual assumed to follow a normal distribution with mean zero and variance $\sigma^2_{e_t}$. In this model parametrization $\beta_t = a_t - d_t$, thus separate estimates of the additive and dominance effects cannot be obtained in this BC population without its reciprocal. Hereafter, we assume an additive model for allelic effects ($d_t = 0$, which leads to $\beta_t = a_t$). We define $\boldsymbol{X} = (X_1, X_2, \cdots, X_n)$, $\boldsymbol{y}_i = (y_{1i}, y_{2i}, \cdots, y_{Ti})'$, $\boldsymbol{Y} = (\boldsymbol{y}_1, \boldsymbol{y}_2, \cdots, \boldsymbol{y}_n)$, $\boldsymbol{e}_i = (e_{1i}, e_{2i}, \cdots, e_{Ti})'$, and $\boldsymbol{E} = (\boldsymbol{e}_1, \boldsymbol{e}_2, \cdots, \boldsymbol{e}_n)$. Then, model (C.15) in matrix form is:

$$\boldsymbol{Y} = \boldsymbol{a}\boldsymbol{X} + \boldsymbol{E} \tag{C.16}$$

The likelihood function of data under model (C.16) is:

$$L(\boldsymbol{a}, \boldsymbol{\Sigma}_e | \boldsymbol{Y}, \boldsymbol{X}, l) = \prod_{i=1}^{n} (2\pi)^{-\frac{T}{2}} |\boldsymbol{\Sigma}_e|^{-\frac{1}{2}} e^{-\frac{1}{2}(\boldsymbol{y}_i - \boldsymbol{a}X_i)' \boldsymbol{\Sigma}_e^{-1}(\boldsymbol{y}_i - \boldsymbol{a}X_i)}$$

The MLE of $\boldsymbol{\Sigma}_e$ and $\boldsymbol{a}$ for the data fitted at a given genetic marker $l$ Morgans

apart from the true position of the QTL are:

$$\hat{\boldsymbol{\Sigma}}_e(l) = \frac{1}{n}(\boldsymbol{Y} - \hat{\boldsymbol{a}}\boldsymbol{X})(\boldsymbol{Y} - \hat{\boldsymbol{a}}\boldsymbol{X})'$$

$$\hat{\boldsymbol{a}}(l) = \boldsymbol{Y}\boldsymbol{X}'(\boldsymbol{X}\boldsymbol{X}')^{-1}$$

$$= \frac{1}{\frac{n}{\sum\limits_{i=1}^{n} X_i^2}} \sum_{i=1}^{n} X_i \boldsymbol{Y}_i$$

$$= \frac{4}{n} \sum_{i=1}^{n} X_i \boldsymbol{Y}_i$$

$$= \bar{\boldsymbol{Y}}_M - \bar{\boldsymbol{Y}}_m$$

The last equality is true under the assumption that the number of subjects with genotypes $MM$ and $Mm$ is equal (i.e., $n_1 = n_2 = n/2$).

The auxiliary equations (equations (C.17), (C.18), (C.19), (C.20) and (C.21)), for obtaining the expected LRT when testing the hypothesis of the QTL located at a genetic marker $l$ Morgans apart from its true position versus the true hypothesis that the QTL is located in position $l = 0$, are given bellow:

$$\sum_{i=1}^{n} (y_{ti} - \hat{a}_t X_i)(y_{t'i} - \hat{a}_{t'} X_i) = \sum_{i=1}^{n} (y_{ti} y_{t'i} - \hat{a}_{t'} X_i y_{ti} - \hat{a}_t X_i y_{t'i} + \hat{a}_t \hat{a}_{t'} X_i^2)$$

$$= \sum_{i=1}^{n} y_{ti} y_{t'i} - \hat{a}_{t'} \sum_{i=1}^{n} X_i y_{ti} - \hat{a}_t \sum_{i=1}^{n} X_i y_{t'i} + \hat{a}_t a_{t'} \sum_{i=1}^{n} X_i^2$$

$$= \sum_{i=1}^{n} y_{ti} y_{t'i} - \frac{n}{4} \hat{a}_{t'} \hat{a}_t - \frac{n}{4} \hat{a}_t \hat{a}_{t'} + \frac{n}{4} \hat{a}_t \hat{a}_{t'}$$

$$= \sum_{i=1}^{n} y_{ti} y_{t'i} - \frac{n}{4} \hat{a}_{t'} \hat{a}_t \qquad (C.17)$$

$$
\begin{aligned}
\mathrm{E}(y_{ti}y_{t'i}) &= \mathrm{E}[(a_t X_i + e_{ti})(a_{t'}X_i + e_{t'i})] \\
&= \mathrm{E}(a_t X_i a_{t'} X_i) + \mathrm{E}(e_{ti}a_{t'}X_i) + \mathrm{E}(e_{t'i}a_t X_i) + \mathrm{E}(e_{t'i}e_{ti}) \\
&= \mathrm{E}(a_t a_{t'} X_i^2) + \mathrm{E}(e_{t'i}e_{ti}) \\
&= \frac{1}{4}a_t a_{t'} + \sigma_{e_{tt'}}
\end{aligned}
\tag{C.18}
$$

$$
\begin{aligned}
\mathrm{E}(\hat{a}_t \hat{a}_{t'}) &= \mathrm{E}[(\bar{y}_{t_M} - \bar{y}_{t_m})(\bar{y}_{t'_M} - \bar{y}_{t'_m})] \\
&= \mathrm{Cov}(\bar{y}_{t_M} - \bar{y}_{t_m}, \bar{y}_{t'_M} - \bar{y}_{t'_m}) + \mathrm{E}(\bar{y}_{t_M} - \bar{y}_{t_m})\mathrm{E}(\bar{y}_{t'_M} - \bar{y}_{t'_m}) \\
&\approx \frac{4}{n}\sigma_{e_{tt'}} + (1 - 2r)^2 a_t a_{t'}
\end{aligned}
\tag{C.19}
$$

From equations (C.17), (C.18) and (C.19) it follows that:

$$
\mathrm{E}[\sum_{i=1}^{n}(y_{ti} - \hat{a}_t X_i)(y_{t'i} - \hat{a}_{t'}X_i)] = (n-1)\sigma_{e_{tt'}} - \frac{n}{4}(1 - e^{-4l})a_t a_{t'}
\tag{C.20}
$$

Assuming $n - 1 \approx n$, which is reasonable under large $n$, equation (C.20) gives the expected values of all elements of the estimated variance-covariance matrix, $\hat{\boldsymbol{\Sigma}}_e(l)$, as follows:

$$
\mathrm{E}[\hat{\boldsymbol{\Sigma}}_e(l)] =
\begin{pmatrix}
\sigma_{e_1}^2 & \sigma_{e_{12}} & \cdots & \sigma_{e_{1T}} \\
\sigma_{e_{21}} & \sigma_{e_2}^2 & \cdots & \sigma_{e_{2T}} \\
\vdots & \vdots & \ddots & \vdots \\
\sigma_{e_{T2}} & \sigma_{e_{T2}} & \cdots & \sigma_{e_T}^2
\end{pmatrix}
+ \frac{1}{4}(1 - e^{-4l})
\begin{pmatrix}
a_1^2 & a_1 a_2 & \cdots & a_1 a_T \\
a_2 a_1 & a_2^2 & \cdots & a_2 a_T \\
\vdots & \vdots & \ddots & \vdots \\
a_T a_1 & a_T a_2 & \cdots & a_T^2
\end{pmatrix}
$$

Assuming $T = 2$, it follows that the determinant of $E[\hat{\boldsymbol{\Sigma}}_e(l)]$ is:

$$
\begin{aligned}
\left|E[\hat{\boldsymbol{\Sigma}}_e(l)]\right| &= [\sigma_{e_1}^2 + \frac{1}{4}(1 - e^{-4l})a_1{}^2][\sigma_{e_2}^2 + \frac{1}{4}(1 - e^{-4l})a_2{}^2] - [\sigma_{e_{12}} + \frac{1}{4}(1 - e^{-4l})a_1 a_2] \\
&= \sigma_{e_1}^2 \sigma_{e_2}^2 (1 - \rho_{12}^2) + \frac{1}{4}(1 - e^{-4l})(\sigma_{e_1}^2 a_2{}^2 - 2\sigma_{e_{12}} a_1 a_2 + \sigma_{e_2}^2 a_1{}^2) \\
&= |\boldsymbol{\Sigma}_e(0)| + |\boldsymbol{\Sigma}_e(0)|\frac{1}{4}(1 - e^{-4l})\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}(0)\boldsymbol{a} \quad\quad (C.21)
\end{aligned}
$$

In fact, equation (C.21) holds good for any number of traits $T$. Therefore, equation (C.21) allows us to derive the approximated expectation of the LRT for testing the hypothesis that the QTL is located at a genetic marker $l$ Morgans apart from its true position versus the true hypothesis that the QTL is located in position $l = 0$. The analytical derivation of the expected value is:

$$
\begin{aligned}
E[\mathrm{LRT}(l)] &= -2E[\hat{\ell}(l) - \hat{\ell}(0)] \\
&= -2E[-\frac{n}{2}\log(|\hat{\boldsymbol{\Sigma}}_e(l)|) - (-\frac{n}{2}\log(|\hat{\boldsymbol{\Sigma}}_e(0)|)] \\
&\approx n(\log E[|\hat{\boldsymbol{\Sigma}}_e(l)|] - \log E[|\hat{\boldsymbol{\Sigma}}_e(0)|]) \\
&= n\log(\frac{E[|\hat{\boldsymbol{\Sigma}}_e(l)|]}{E[|\hat{\boldsymbol{\Sigma}}_e(0)|]}) \\
&\approx n\log(\frac{|E[\hat{\boldsymbol{\Sigma}}_e(l)]|}{|E[\hat{\boldsymbol{\Sigma}}_e(0)]|}) \\
&= n\log[1 + \frac{1}{4}(1 - e^{-4l})\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}(0)\boldsymbol{a}] \\
&\approx \frac{n}{4}(1 - e^{-4l})\boldsymbol{a}'\boldsymbol{\Sigma}_e^{-1}(0)\boldsymbol{a}
\end{aligned}
$$

# D

# Evaluation of the MTMIM model by simulation study and experimental data analysis

## D.1 Complementary tables

Table D.1: Decomposition of total power of QTL identification in scenario SII (Table 3.10) into QTL-by-trait power for 10% genome-wide significance level. Subsets (1, 0, 0), (1, 1, 0) and (1, 1, 1) contain replicates with QTL affecting T1 only, T1 and T2, and T1, T2 and T3, respectively. We show only three subsets out of 7 subsets of the full decomposition. These three subsets account for most of the space, besides they are the most interesting in scenario SII.

| Subsets | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| | Q1 | 42.4 | 45.4 | 46.4 | 56.8 | 60.2 | 61.2 | 63.0 | 66.4 | 67.6 |
| | Q2 | 0.4 | 0.6 | 0.6 | 0.8 | 1.0 | 1.0 | 1.0 | 1.2 | 1.2 |
| (1,0,0) | Q3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| | Q4 | 0.2 | 0.2 | 0.2 | 0.6 | 0.6 | 0.6 | 0.8 | 0.8 | 0.8 |
| | Q5 | 44.2 | 48.0 | 48.6 | 55.4 | 59.6 | 60.2 | 59.8 | 64.0 | 64.8 |
| | Q1 | 2.8 | 2.8 | 2.8 | 3.8 | 3.8 | 3.8 | 4.2 | 4.2 | 4.2 |
| | Q2 | 76.4 | 80.6 | 81.4 | 81.8 | 85.6 | 86.6 | 82.8 | 86.4 | 87.4 |
| (1,1,0) | Q3 | 3.6 | 3.8 | 4.0 | 4.2 | 4.6 | 4.8 | 4.6 | 5.0 | 5.2 |
| | Q4 | 79.2 | 81.6 | 82.6 | 83.4 | 86.2 | 87.0 | 84.0 | 87.2 | 88.0 |
| | Q5 | 4.0 | 4.8 | 5.2 | 5.6 | 6.6 | 7.0 | 6.8 | 8.2 | 8.6 |
| | Q1 | 0.6 | 0.6 | 0.6 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 | 0.8 |
| | Q2 | 6.2 | 6.8 | 6.8 | 6.0 | 6.4 | 6.6 | 6.0 | 6.6 | 6.8 |
| (1,1,1) | Q3 | 87.6 | 91.0 | 93.0 | 86.0 | 89.8 | 91.8 | 85.2 | 89.0 | 90.8 |
| | Q4 | 5.4 | 5.6 | 6.0 | 5.8 | 6.0 | 6.4 | 5.8 | 5.8 | 6.4 |
| | Q5 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.2: Decomposition of total power of QTL identification in scenario SIII (Table 3.11) into QTL-by-trait power for 10% genome-wide significance level. Subsets (1, 0), (0, 1) and (1, 1) contain replicates with QTL affecting T1 only, T2 only, and T1 and T2, respectively.

| Subsets | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| (1,0) | Q1 | 33.6 | 34.6 | 34.6 | 33.4 | 34.0 | 34.2 | 36.2 | 36.8 | 37.0 |
| | Q2 | 1.4 | 1.4 | 1.6 | 5.0 | 5.0 | 5.2 | 2.8 | 2.8 | 3.0 |
| | Q3 | 3.0 | 2.4 | 2.6 | 3.2 | 2.8 | 2.8 | 3.8 | 3.4 | 3.2 |
| | Q4 | 0.6 | 0.8 | 1.0 | 1.0 | 1.0 | 1.2 | 1.2 | 1.0 | 1.2 |
| | Q5 | 46.6 | 46.6 | 47.0 | 46.2 | 46.8 | 47.2 | 45.0 | 46.0 | 46.6 |
| (0,1) | Q1 | 1.8 | 1.6 | 1.6 | 2.4 | 2.6 | 2.6 | 2.4 | 2.8 | 2.8 |
| | Q2 | 33.2 | 34.0 | 34.0 | 32.2 | 33.0 | 33.2 | 35.4 | 36.2 | 36.8 |
| | Q3 | 3.2 | 3.8 | 3.8 | 3.4 | 3.8 | 4.0 | 3.6 | 4.0 | 4.4 |
| | Q4 | 48.6 | 48.6 | 49.2 | 48.4 | 48.6 | 49.6 | 48.8 | 49.6 | 50.6 |
| | Q5 | 1.2 | 1.6 | 1.2 | 1.2 | 1.2 | 0.8 | 1.6 | 1.2 | 0.8 |
| (1,1) | Q1 | 28.0 | 29.2 | 29.4 | 27.6 | 28.6 | 29.0 | 29.2 | 30.4 | 30.8 |
| | Q2 | 28.8 | 29.2 | 29.4 | 28.2 | 28.6 | 28.8 | 28.4 | 29.0 | 29.2 |
| | Q3 | 84.8 | 88.2 | 89.6 | 86.2 | 89.8 | 91.2 | 86.2 | 89.6 | 91.2 |
| | Q4 | 23.4 | 25.4 | 25.4 | 26.2 | 27.8 | 28.0 | 26.0 | 27.6 | 27.8 |
| | Q5 | 16.6 | 17.4 | 17.4 | 17.4 | 18.2 | 18.6 | 19.8 | 20.8 | 21.0 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.3: Mean of QTL positions (cM) in MIM and MTMIM models in scenario SI across significance levels (1, 5 and 10%). Range of the standard error (SE) of means is also shown.

| Analysis (trait) | QTL | Position | 1% | | | 5% | | | 10% | | | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 | |
| MIM (T1) | Q1 | 23 | 23.5 | 23.9 | 24.0 | 23.5 | 23.7 | 23.8 | 23.4 | 23.6 | 23.6 | 0.27-0.39 |
| | Q2 | 15 | 14.4 | 14.8 | 14.9 | 14.5 | 14.8 | 14.9 | 15.7 | 14.9 | 15.3 | 0.31-0.38 |
| | Q3 | 45 | 45.5 | 45.2 | 45.6 | 45.5 | 45.2 | 45.6 | 45.6 | 45.3 | 45.6 | 0.32-0.42 |
| | Q4 | 67 | 67.2 | 66.8 | 66.8 | 66.9 | 66.6 | 66.6 | 66.8 | 66.5 | 66.5 | 0.26-0.23 |
| | Q5 | 53 | 52.9 | 52.8 | 52.8 | 52.7 | 52.5 | 52.4 | 52.6 | 52.4 | 52.5 | 0.30-0.36 |
| MIM (T2) | Q1 | 23 | 23.4 | 23.5 | 23.4 | 23.6 | 23.8 | 23.8 | 23.7 | 23.9 | 23.9 | 0.32-0.39 |
| | Q2 | 15 | 13.4 | 13.5 | 13.5 | 14.1 | 14.1 | 14.4 | 14.1 | 14.1 | 14.4 | 0.29-0.36 |
| | Q3 | 45 | 45.0 | 45.2 | 45.2 | 44.9 | 45.1 | 45.0 | 44.8 | 44.9 | 45.0 | 0.29-0.39 |
| | Q4 | 67 | 66.8 | 66.5 | 66.2 | 66.7 | 66.4 | 66.2 | 66.8 | 66.3 | 66.2 | 0.27-0.39 |
| | Q5 | 53 | 52.3 | 52.4 | 52.4 | 52.6 | 52.5 | 52.4 | 52.4 | 52.4 | 52.4 | 0.27-0.38 |
| MIM (T3) | Q1 | 23 | 23.6 | 23.6 | 23.9 | 23.7 | 23.8 | 24.1 | 23.7 | 24.0 | 24.3 | 0.29-0.41 |
| | Q2 | 15 | 14.0 | 14.3 | 14.7 | 13.9 | 14.2 | 14.6 | 14.2 | 14.6 | 14.9 | 0.29-0.39 |
| | Q3 | 45 | 44.9 | 45.0 | 45.2 | 44.8 | 44.9 | 45.2 | 44.7 | 44.8 | 45.1 | 0.30-0.43 |
| | Q4 | 67 | 67.0 | 66.8 | 66.5 | 67.1 | 66.9 | 66.7 | 67.3 | 67.1 | 66.8 | 0.28-0.39 |
| | Q5 | 53 | 52.8 | 52.7 | 52.6 | 52.5 | 52.3 | 52.3 | 52.6 | 52.4 | 52.3 | 0.31-0.42 |
| MTMIM (T1,T2,T3) | Q1 | 23 | 23.61 | 23.67 | 23.67 | 23.59 | 23.73 | 23.74 | 23.62 | 23.76 | 23.75 | 0.13-0.17 |
| | Q2 | 15 | 14.13 | 14.18 | 14.20 | 14.15 | 14.21 | 14.22 | 14.15 | 14.20 | 14.22 | 0.15-0.16 |
| | Q3 | 45 | 45.47 | 45.54 | 45.65 | 45.48 | 45.55 | 45.66 | 45.50 | 45.57 | 45.68 | 0.16-0.18 |
| | Q4 | 67 | 67.69 | 67.60 | 67.57 | 67.71 | 67.61 | 67.58 | 67.68 | 67.66 | 67.62 | 0.14-0.17 |
| | Q5 | 53 | 52.85 | 52.80 | 52.85 | 52.82 | 52.78 | 52.83 | 52.76 | 52.73 | 52.82 | 0.14-0.16 |

[a] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.4: Mean of positions of QTL (cM) in MIM and MTMIM models in scenario SII across significance levels (1, 5 and 10%). Range of the standard error (SE) of means is also shown.

| Analysis (trait) | QTL | Position | 1% | | | 5% | | | 10% | | | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 | |
| MIM (T1) | Q1 | 23 | 23.4 | 23.4 | 23.4 | 23.3 | 23.5 | 23.6 | 23.4 | 23.7 | 23.8 | 0.27-0.36 |
| | Q2 | 15 | 14.4 | 14.4 | 14.6 | 14.3 | 14.5 | 14.6 | 14.4 | 14.6 | 14.9 | 0.29-0.38 |
| | Q3 | 45 | 45.5 | 45.3 | 45.4 | 45.5 | 45.3 | 45.3 | 45.4 | 45.4 | 45.2 | 0.32-0.42 |
| | Q4 | 67 | 67.2 | 66.9 | 66.8 | 67.1 | 66.9 | 66.6 | 67.1 | 66.9 | 66.7 | 0.27-0.34 |
| | Q5 | 53 | 52.9 | 52.7 | 52.5 | 53.0 | 52.9 | 52.7 | 52.9 | 52.9 | 52.8 | 0.28-0.39 |
| MIM (T2) | Q2 | 15 | 14.3 | 14.5 | 14.7 | 14.3 | 14.5 | 14.7 | 14.5 | 14.7 | 14.9 | 0.27-0.34 |
| | Q3 | 45 | 45.4 | 45.3 | 45.4 | 45.4 | 45.1 | 45.3 | 45.6 | 45.2 | 45.4 | 0.30-0.39 |
| | Q4 | 67 | 67.3 | 66.9 | 66.7 | 67.4 | 67.1 | 66.8 | 67.2 | 67.0 | 66.6 | 0.25-0.35 |
| MIM (T3) | Q3 | 45 | 44.9 | 44.8 | 44.9 | 44.7 | 44.9 | 45.0 | 44.7 | 44.7 | 44.8 | 0.35-0.47 |
| MTMIM (T1,T2,T3) | Q1 | 23 | 23.4 | 23.2 | 23.1 | 23.5 | 23.4 | 23.3 | 23.4 | 23.5 | 23.5 | 0.31-0.37 |
| | Q2 | 15 | 14.3 | 14.4 | 14.5 | 14.4 | 14.4 | 14.5 | 14.4 | 14.4 | 14.5 | 0.21-0.23 |
| | Q3 | 45 | 45.1 | 45.1 | 45.1 | 45.0 | 45.0 | 45.0 | 44.9 | 44.9 | 44.9 | 0.18-0.19 |
| | Q4 | 67 | 67.6 | 67.5 | 67.4 | 67.6 | 67.5 | 67.5 | 67.6 | 67.6 | 67.5 | 0.18-0.21 |
| | Q5 | 53 | 52.9 | 52.7 | 52.8 | 52.8 | 52.6 | 52.8 | 52.9 | 52.8 | 52.9 | 0.33-0.42 |

[a] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.5: Mean of positions of QTL (cM) in MIM and MTMIM models in scenario SIII across significance levels (1, 5 and 10%). Range of the standard error (SE) of means is also shown.

| Analysis (trait) | QTL | Position | 1% | | | 5% | | | 10% | | | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 | |
| MIM (T1) | Q1 | 23 | 22.3 | 22.1 | 22.1 | 22.1 | 21.9 | 21.9 | 22.1 | 22.1 | 22.0 | 0.18-0.23 |
| | Q3 | 45 | 44.2 | 44.2 | 44.2 | 44.5 | 44.4 | 44.2 | 44.3 | 44.3 | 44.1 | 0.27-0.35 |
| | Q5 | 53 | 52.5 | 53.1 | 53.1 | 52.3 | 52.8 | 52.8 | 52.4 | 52.9 | 52.9 | 0.25-0.28 |
| MIM (T2) | Q2 | 33 | 34.9 | 35.3 | 35.5 | 34.9 | 35.6 | 35.7 | 34.9 | 35.5 | 35.6 | 0.21-0.30 |
| | Q3 | 45 | 43.9 | 44.1 | 43.9 | 43.7 | 43.8 | 43.8 | 43.8 | 43.8 | 43.8 | 0.30-0.37 |
| | Q4 | 38 | 37.1 | 36.8 | 36.4 | 36.9 | 36.4 | 36.2 | 36.9 | 36.5 | 36.2 | 0.21-0.30 |
| MTMIM (T1,T2) | Q1 | 23 | 23.1 | 22.9 | 22.9 | 23.7 | 23.6 | 23.5 | 23.2 | 23.1 | 23.0 | 0.20-0.25 |
| | Q2 | 33 | 33.4 | 33.5 | 33.5 | 33.7 | 33.7 | 33.7 | 33.5 | 33.6 | 33.7 | 0.21-0.27 |
| | Q3 | 45 | 44.5 | 44.6 | 44.6 | 44.5 | 44.6 | 44.5 | 44.4 | 44.6 | 44.5 | 0.19-0.23 |
| | Q4 | 38 | 38.8 | 38.4 | 38.4 | 38.8 | 38.6 | 38.6 | 38.6 | 38.5 | 38.5 | 0.21-0.27 |
| | Q5 | 53 | 51.3 | 51.3 | 51.4 | 51.3 | 51.5 | 51.7 | 51.4 | 51.7 | 51.8 | 0.26-0.29 |

[a] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.6: Mean length (cM) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SI across significance levels (1, 5 and 10%). Range of the standard error (SE) of means is also shown.

| Analysis (trait) | QTL | 1% | | | 5% | | | 10% | | | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 | |
| MIM (T1) | Q1 | 18.9 | 25.7 | 32.4 | 19.9 | 27.2 | 34.7 | 20.7 | 28.0 | 35.7 | 0.35-0.64 |
| | Q2 | 19.2 | 25.1 | 30.5 | 20.6 | 27.3 | 34.0 | 21.3 | 28.4 | 35.5 | 0.37-0.74 |
| | Q3 | 19.9 | 27.9 | 35.4 | 21.3 | 29.5 | 38.0 | 22.3 | 31.1 | 39.8 | 0.42-0.78 |
| | Q4 | 18.2 | 23.7 | 29.0 | 19.4 | 25.7 | 31.8 | 19.9 | 26.4 | 33.4 | 0.33-0.74 |
| | Q5 | 19.8 | 27.1 | 33.9 | 20.9 | 28.7 | 36.3 | 21.3 | 29.1 | 36.4 | 0.40-0.63 |
| MIM (T2) | Q1 | 20.0 | 26.9 | 34.2 | 21.3 | 29.0 | 36.4 | 22.6 | 30.7 | 38.4 | 0.40-0.69 |
| | Q2 | 19.5 | 25.9 | 31.6 | 20.9 | 27.6 | 34.2 | 21.6 | 28.4 | 35.5 | 0.39-0.74 |
| | Q3 | 20.2 | 27.4 | 35.3 | 21.3 | 29.6 | 38.5 | 21.9 | 30.4 | 39.2 | 0.39-0.77 |
| | Q4 | 19.1 | 24.6 | 30.0 | 19.9 | 25.9 | 32.1 | 20.2 | 26.6 | 33.2 | 0.38-0.67 |
| | Q5 | 17.5 | 26.3 | 33.2 | 20.4 | 28.1 | 35.4 | 21.6 | 29.6 | 36.7 | 0.37-0.66 |
| MIM (T3) | Q1 | 19.9 | 26.6 | 33.2 | 21.2 | 28.8 | 36.0 | 21.7 | 29.9 | 37.2 | 0.41-0.68 |
| | Q2 | 19.2 | 24.6 | 30.2 | 20.5 | 26.5 | 33.4 | 21.3 | 27.7 | 34.9 | 0.36-0.70 |
| | Q3 | 20.9 | 28.8 | 36.8 | 22.1 | 30.5 | 39.4 | 22.4 | 31.2 | 40.3 | 0.43-0.76 |
| | Q4 | 18.6 | 23.8 | 29.3 | 19.3 | 25.1 | 31.5 | 19.6 | 27.8 | 33.3 | 0.35-0.69 |
| | Q5 | 20.2 | 26.5 | 32.8 | 21.5 | 28.8 | 35.5 | 21.9 | 29.5 | 36.4 | 0.41-0.63 |
| MTMIM (T1,T2,T3) | Q1 | 12.38 | 15.88 | 19.00 | 12.42 | 16.04 | 19.13 | 12.40 | 16.01 | 19.11 | 0.17-0.30 |
| | Q2 | 12.60 | 16.11 | 19.54 | 12.56 | 16.04 | 19.47 | 12.55 | 16.02 | 19.45 | 0.17-0.27 |
| | Q3 | 12.41 | 15.90 | 19.31 | 12.44 | 15.90 | 19.36 | 12.51 | 15.98 | 19.44 | 0.15-0.28 |
| | Q4 | 12.09 | 15.31 | 18.09 | 12.13 | 15.35 | 18.18 | 12.21 | 15.34 | 18.19 | 0.15-0.23 |
| | Q5 | 12.32 | 15.66 | 18.73 | 12.30 | 15.62 | 18.17 | 12.25 | 15.58 | 18.65 | 0.15-0.25 |

[a] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.7: Mean length (cM) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SII across significance levels (1, 5 and 10%). Range of the standard error (SE) of means is also shown.

| Analysis (trait) | QTL | 1% | | | 5% | | | 10% | | | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 | |
| MIM (T1) | Q1 | 19.9 | 27.0 | 33.9 | 21.2 | 28.5 | 36.2 | 21.7 | 29.4 | 37.3 | 0.38-0.66 |
| | Q2 | 19.3 | 25.1 | 30.4 | 20.6 | 26.8 | 33.3 | 21.1 | 27.7 | 34.9 | 0.37-0.73 |
| | Q3 | 20.9 | 28.5 | 36.3 | 23.1 | 31.9 | 40.4 | 23.7 | 33.0 | 41.9 | 0.42-0.81 |
| | Q4 | 18.7 | 24.2 | 30.3 | 19.8 | 26.1 | 33.6 | 20.2 | 26.7 | 35.4 | 0.35-0.79 |
| | Q5 | 18.9 | 25.9 | 32.8 | 20.6 | 27.9 | 35.6 | 21.3 | 28.7 | 36.4 | 0.37-0.68 |
| MIM (T2) | Q2 | 19.6 | 25.5 | 30.7 | 20.6 | 26.9 | 32.7 | 21.0 | 27.9 | 34.1 | 0.35-0.67 |
| | Q3 | 20.9 | 27.8 | 35.4 | 21.8 | 29.2 | 37.7 | 22.3 | 29.8 | 39.1 | 0.39-0.64 |
| | Q4 | 18.3 | 24.0 | 29.7 | 19.2 | 25.3 | 31.3 | 19.6 | 26.1 | 32.6 | 0.30-0.67 |
| MIM (T3) | Q3 | 22.2 | 29.8 | 38.9 | 24.5 | 33.8 | 44.1 | 25.3 | 35.3 | 46.2 | 0.53-0.88 |
| MTMIM (T1, T2, T3) | Q1 | 17.9 | 23.6 | 29.6 | 19.2 | 25.5 | 32.0 | 20.0 | 26.4 | 33.1 | 0.36-0.57 |
| | Q2 | 15.9 | 20.3 | 24.7 | 16.1 | 20.9 | 25.2 | 16.2 | 21.0 | 25.3 | 0.24-0.39 |
| | Q3 | 13.1 | 17.1 | 20.7 | 13.1 | 17.1 | 20.7 | 13.1 | 17.2 | 20.7 | 0.18-0.33 |
| | Q4 | 15.4 | 19.9 | 23.7 | 15.6 | 20.3 | 23.9 | 15.6 | 20.3 | 24.2 | 0.23-0.39 |
| | Q5 | 17.7 | 23.5 | 29.1 | 18.9 | 25.1 | 31.4 | 19.7 | 26.1 | 32.6 | 0.36-0.60 |

[a] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.8: Mean length (cM) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SIII across significance levels (1, 5 and 10%). Range of the standard error (SE) of means is also shown.

| Analysis (trait) | QTL | 1% | | | 5% | | | 10% | | | SE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[a] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 | |
| MIM (T1) | Q1 | 18.8 | 24.7 | 31.6 | 19.8 | 26.1 | 33.9 | 20.2 | 26.8 | 34.9 | 0.34-0.74 |
| | Q3 | 19.7 | 26.7 | 33.9 | 20.7 | 28.5 | 36.8 | 21.1 | 29.1 | 37.8 | 0.37-0.74 |
| | Q5 | 19.1 | 25.8 | 32.9 | 19.8 | 27.0 | 34.8 | 20.4 | 27.8 | 35.9 | 0.35-0.75 |
| MIM (T2) | Q2 | 19.2 | 25.9 | 33.3 | 20.3 | 28.3 | 36.5 | 20.6 | 28.9 | 37.5 | 0.40-0.89 |
| | Q3 | 20.0 | 26.8 | 34.2 | 21.8 | 29.0 | 37.5 | 22.1 | 29.6 | 38.4 | 0.39-0.77 |
| | Q4 | 18.1 | 25.0 | 32.5 | 19.2 | 26.4 | 34.9 | 19.5 | 26.9 | 35.8 | 0.35-0.75 |
| MTMIM (T1,T2) | Q1 | 26.7 | 36.2 | 46.0 | 28.1 | 37.9 | 47.8 | 28.8 | 38.4 | 47.6 | 1.10-1.45 |
| | Q2 | 26.4 | 37.1 | 46.4 | 28.4 | 39.4 | 48.7 | 28.1 | 39.4 | 48.7 | 1.12-1.44 |
| | Q3 | 16.7 | 22.5 | 28.3 | 16.9 | 22.8 | 28.7 | 16.8 | 22.9 | 28.8 | 0.32-0.65 |
| | Q4 | 26.8 | 38.2 | 48.9 | 26.6 | 37.8 | 48.2 | 26.7 | 37.3 | 47.6 | 0.92-1.26 |
| | Q5 | 29.7 | 41.1 | 51.3 | 29.5 | 40.5 | 50.4 | 28.9 | 39.9 | 49.5 | 1.06-1.37 |

[a] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.9: Coverage (%) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SI.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 91.9 | 97.1 | 99.4 | 92.9 | 97.4 | 99.3 | 92.9 | 97.9 | 100 |
| | Q2 | 91.0 | 95.5 | 96.4 | 91.5 | 96.2 | 97.2 | 92.3 | 96.3 | 97.6 |
| | Q3 | 90.0 | 96.6 | 98.6 | 91.5 | 96.7 | 98.8 | 90.9 | 96.7 | 98.5 |
| | Q4 | 91.8 | 96.8 | 99.1 | 92.7 | 96.7 | 98.8 | 92.5 | 96.5 | 98.5 |
| | Q5 | 90.6 | 95.2 | 96.9 | 91.5 | 96.5 | 98.2 | 91.3 | 96.4 | 98.2 |
| MIM (T2) | Q1 | 94.9 | 97.9 | 98.8 | 94.8 | 99.3 | 100 | 95.9 | 99.6 | 100 |
| | Q2 | 93.4 | 96.7 | 97.9 | 93.7 | 96.9 | 98.3 | 93.4 | 96.6 | 98.2 |
| | Q3 | 91.4 | 94.3 | 97.9 | 91.7 | 95.0 | 98.3 | 91.8 | 94.9 | 98.2 |
| | Q4 | 90.9 | 97.1 | 97.9 | 91.0 | 97.6 | 98.3 | 90.8 | 97.3 | 98.0 |
| | Q5 | 91.2 | 96.6 | 98.9 | 90.7 | 96.1 | 99.3 | 92.2 | 96.7 | 99.4 |
| MIM (T3) | Q1 | 87.1 | 96.2 | 98.2 | 89.3 | 96.5 | 98.8 | 90.3 | 96.9 | 99.1 |
| | Q2 | 93.2 | 97.5 | 99.1 | 93.7 | 97.3 | 98.3 | 94.1 | 97.5 | 98.7 |
| | Q3 | 86.7 | 95.1 | 97.4 | 87.2 | 95.3 | 97.6 | 87.4 | 94.6 | 96.9 |
| | Q4 | 89.5 | 95.6 | 97.9 | 92.1 | 96.7 | 98.3 | 92.7 | 96.9 | 98.5 |
| | Q5 | 90.1 | 97.3 | 99.4 | 94.1 | 98.1 | 100 | 93.8 | 97.8 | 100 |
| MTMIM (T1,T2,T3) | Q1 | 96.6 | 99.2 | 99.4 | 96.4 | 99.4 | 99.6 | 96.4 | 99.4 | 99.6 |
| | Q2 | 97.0 | 98.0 | 98.6 | 97.0 | 98.0 | 98.6 | 96.8 | 98.2 | 98.8 |
| | Q3 | 94.8 | 97.6 | 98.9 | 94.8 | 97.6 | 99.0 | 94.6 | 97.4 | 98.8 |
| | Q4 | 94.0 | 98.8 | 99.4 | 94.0 | 98.8 | 99.4 | 93.8 | 99.0 | 99.4 |
| | Q5 | 96.4 | 98.6 | 99.6 | 96.4 | 98.6 | 99.6 | 96.0 | 98.6 | 99.6 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.10: Coverage (%) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SII.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 90.1 | 95.6 | 98.6 | 91.1 | 95.9 | 99.3 | 91.4 | 95.7 | 99.3 |
| | Q2 | 91.9 | 96.5 | 98.4 | 92.3 | 96.8 | 98.0 | 92.2 | 95.8 | 98.1 |
| | Q3 | 88.8 | 95.5 | 97.3 | 89.6 | 96.2 | 98.1 | 88.8 | 95.8 | 98.2 |
| | Q4 | 92.7 | 96.9 | 98.8 | 92.2 | 95.6 | 98.5 | 92.2 | 95.8 | 98.4 |
| | Q5 | 93.1 | 98.5 | 99.7 | 93.3 | 98.6 | 100 | 93.4 | 98.8 | 99.6 |
| MIM (T2) | Q2 | 92.9 | 97.9 | 98.9 | 92.5 | 97.5 | 98.6 | 92.6 | 97.4 | 98.7 |
| | Q3 | 90.3 | 95.5 | 98.0 | 90.4 | 95.8 | 98.2 | 90.6 | 95.9 | 98.3 |
| | Q4 | 93.9 | 97.7 | 99.8 | 94.6 | 97.8 | 99.6 | 95.3 | 98.1 | 99.6 |
| MIM (T3) | Q3 | 89.1 | 94.0 | 96.8 | 89.6 | 94.1 | 96.5 | 88.8 | 94.6 | 96.8 |
| MTMIM (T1,T2,T3) | Q1 | 87.4 | 94.4 | 96.8 | 88.9 | 95.4 | 97.3 | 89.5 | 95.6 | 97.6 |
| | Q2 | 92.2 | 98.0 | 98.9 | 93.0 | 97.7 | 98.9 | 93.1 | 97.8 | 98.9 |
| | Q3 | 93.4 | 97.0 | 99.4 | 92.8 | 97.0 | 99.4 | 92.8 | 97.2 | 99.4 |
| | Q4 | 94.4 | 97.3 | 98.9 | 94.6 | 97.7 | 98.9 | 94.2 | 97.5 | 98.9 |
| | Q5 | 88.4 | 97.3 | 99.7 | 89.4 | 97.8 | 99.7 | 89.5 | 97.8 | 99.8 |

[a] Genome-wide significance level.
[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.11: Coverage (%) of LOD-$d$ support interval for position of QTL in MIM and MTMIM models in scenario SIII.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 97.9 | 99.4 | 99.7 | 97.9 | 99.5 | 99.7 | 98.0 | 99.8 | 100 |
| | Q3 | 92.1 | 95.9 | 97.9 | 92.1 | 96.0 | 98.2 | 92.1 | 96.2 | 98.5 |
| | Q5 | 95.2 | 99.4 | 99.4 | 95.2 | 99.2 | 99.2 | 95.4 | 99.3 | 99.3 |
| MIM (T2) | Q2 | 96.0 | 98.5 | 99.4 | 95.2 | 98.7 | 99.2 | 94.9 | 98.5 | 98.9 |
| | Q3 | 92.2 | 96.5 | 98.5 | 92.8 | 96.3 | 98.3 | 92.4 | 96.2 | 98.1 |
| | Q4 | 95.1 | 96.9 | 98.9 | 95.5 | 97.7 | 99.3 | 95.6 | 97.8 | 99.3 |
| MTMIM (T1,T2) | Q1 | 97.8 | 100 | 100 | 98.8 | 100 | 100 | 97.9 | 100 | 100 |
| | Q2 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| | Q3 | 93.8 | 97.32 | 98.9 | 93.7 | 97.4 | 98.9 | 93.9 | 97.4 | 99.2 |
| | Q4 | 98.9 | 100 | 100 | 98.4 | 100 | 100 | 97.4 | 100 | 100 |
| | Q5 | 100 | 100 | 100 | 97.6 | 99.7 | 100 | 97.1 | 99.4 | 100 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

Table D.12: Mean of effect of QTL in MIM and MTMIM models in scenario SI.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 0.60[c] | 0.60 | 0.60 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 |
| | Q2 | 0.60 | 0.60 | 0.60 | 0.57 | 0.57 | 0.57 | 0.56 | 0.56 | 0.56 |
| | Q3 | 0.59 | 0.59 | 0.60 | 0.57 | 0.57 | 0.57 | 0.56 | 0.56 | 0.55 |
| | Q4 | 0.58 | 0.59 | 0.59 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 |
| | Q5 | 0.59 | 0.59 | 0.59 | 0.56 | 0.56 | 0.56 | 0.55 | 0.56 | 0.56 |
| MIM (T2) | Q1 | 0.59 | 0.59 | 0.59 | 0.57 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 |
| | Q2 | 0.59 | 0.59 | 0.59 | 0.56 | 0.56 | 0.56 | 0.55 | 0.56 | 0.56 |
| | Q3 | 0.60 | 0.60 | 0.60 | 0.57 | 0.57 | 0.57 | 0.56 | 0.56 | 0.56 |
| | Q4 | 0.59 | 0.59 | 0.59 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 |
| | Q5 | 0.59 | 0.59 | 0.59 | 0.57 | 0.57 | 0.57 | 0.56 | 0.55 | 0.56 |
| MIM (T3) | Q1 | 0.59 | 0.59 | 0.59 | 0.56 | 0.57 | 0.57 | 0.56 | 0.56 | 0.56 |
| | Q2 | 0.60 | 0.60 | 0.60 | 0.56 | 0.57 | 0.57 | 0.55 | 0.55 | 0.55 |
| | Q3 | 0.59 | 0.59 | 0.59 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 |
| | Q4 | 0.59 | 0.59 | 0.59 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 |
| | Q5 | 0.60 | 0.60 | 0.60 | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 |
| MTMIM (T1) | Q1 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.51 | 0.51 |
| | Q2 | 0.51 | 0.52 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 |
| | Q3 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 |
| | Q4 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 |
| | Q5 | 0.52 | 0.52 | 0.52 | 0.51 | 0.51 | 0.51 | 0.52 | 0.52 | 0.51 |
| MTMIM (T2) | Q1 | 0.51 | 0.51 | 0.51 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| | Q2 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 |
| | Q3 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 |
| | Q4 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.50 | 0.50 |
| | Q5 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 |
| MTMIM (T3) | Q1 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 |
| | Q2 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 |
| | Q3 | 0.51 | 0.52 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 |
| | Q4 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.52 | 0.51 |
| | Q5 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.53 | 0.53 | 0.52 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

[c] The standard error of the means in this table ranges from 0.005 to 0.008.

Table D.13: Mean of effect of QTL in MIM and MTMIM models in scenario SII.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 0.59[c] | 0.59 | 0.59 | 0.57 | 0.57 | 0.57 | 0.56 | 0.56 | 0.56 |
| | Q2 | 0.60 | 0.60 | 0.60 | 0.56 | 0.56 | 0.56 | 0.55 | 0.56 | 0.56 |
| | Q3 | 0.58 | 0.58 | 0.58 | 0.55 | 0.55 | 0.55 | 0.54 | 0.54 | 0.54 |
| | Q4 | 0.60 | 0.60 | 0.60 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 |
| | Q5 | 0.59 | 0.59 | 0.60 | 0.56 | 0.56 | 0.57 | 0.55 | 0.55 | 0.56 |
| MIM (T2) | Q2 | 0.59 | 0.59 | 0.59 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 |
| | Q3 | 0.59 | 0.59 | 0.59 | 0.57 | 0.58 | 0.58 | 0.57 | 0.57 | 0.57 |
| | Q4 | 0.60 | 0.60 | 0.60 | 0.58 | 0.58 | 0.58 | 0.58 | 0.57 | 0.57 |
| MIM (T3) | Q3 | 0.54 | 0.54 | 0.54 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 |
| MTMIM (T1) | Q1 | 0.60 | 0.60 | 0.60 | 0.58 | 0.58 | 0.58 | 0.56 | 0.56 | 0.56 |
| | Q2 | 0.53 | 0.53 | 0.53 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 |
| | Q3 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.50 |
| | Q4 | 0.54 | 0.54 | 0.54 | 0.53 | 0.53 | 0.53 | 0.52 | 0.52 | 0.52 |
| | Q5 | 0.60 | 0.60 | 0.60 | 0.58 | 0.58 | 0.58 | 0.56 | 0.56 | 0.56 |
| MTMIM (T2) | Q1 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | Q2 | 0.56 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 | 0.54 | 0.54 |
| | Q3 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 |
| | Q4 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 |
| | Q5 | 0.00 | 0.01 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| MTMIM (T3) | Q1 | 0.01 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | Q2 | 0.00 | 0.01 | 0.01 | 0.00 | 0.01 | 0.01 | 0.00 | 0.01 | 0.01 |
| | Q3 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.43 |
| | Q4 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | Q5 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

[c] The standard error of the means in this table ranges from 0.003 to 0.007.

Table D.14: Mean of effect of QTL in MIM and MTMIM models in scenario SIII.

| Analysis (trait) | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| MIM (T1) | Q1 | 0.60[c] | 0.60 | 0.60 | 0.58 | 0.58 | 0.58 | 0.57 | 0.57 | 0.57 |
| | Q3 | 0.60 | 0.60 | 0.60 | 0.58 | 0.58 | 0.58 | 0.57 | 0.57 | 0.57 |
| | Q5 | 0.60 | 0.60 | 0.60 | 0.58 | 0.58 | 0.58 | 0.58 | 0.58 | 0.58 |
| MIM (T2) | Q2 | 0.61 | 0.61 | 0.61 | 0.59 | 0.59 | 0.59 | 0.58 | 0.58 | 0.58 |
| | Q3 | 0.60 | 0.60 | 0.60 | 0.57 | 0.58 | 0.57 | 0.57 | 0.57 | 0.57 |
| | Q4 | 0.61 | 0.61 | 0.61 | 0.58 | 0.58 | 0.58 | 0.58 | 0.58 | 0.58 |
| MTMIM (T1) | Q1 | 0.56 | 0.57 | 0.57 | 0.55 | 0.55 | 0.55 | 0.57 | 0.56 | 0.57 |
| | Q2 | 0.24 | 0.23 | 0.24 | 0.26 | 0.25 | 0.25 | 0.21 | 0.20 | 0.20 |
| | Q3 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.52 | 0.52 |
| | Q4 | 0.14 | 0.15 | 0.15 | 0.13 | 0.14 | 0.13 | 0.11 | 0.13 | 0.12 |
| | Q5 | 0.58 | 0.58 | 0.58 | 0.58 | 0.59 | 0.59 | 0.58 | 0.58 | 0.59 |
| MTMIM (T2) | Q1 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.24 | 0.23 | 0.24 |
| | Q2 | 0.57 | 0.57 | 0.57 | 0.53 | 0.53 | 0.53 | 0.55 | 0.55 | 0.55 |
| | Q3 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 | 0.54 | 0.53 | 0.54 | 0.54 |
| | Q4 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.58 | 0.59 | 0.60 | 0.58 |
| | Q5 | 0.11 | 0.11 | 0.11 | 0.10 | 0.10 | 0.11 | 0.09 | 0.09 | 0.09 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

[c] The standard error of the means in this table ranges from 0.007 to 0.019.

Table D.15: Frequency of rejection of pleiotropy model for each pleiotropic QTL in MTMIM model in scenario SIII. The AICc and LRT criteria are compared.

| Criterion | QTL | 1%[a] | | | 5% | | | 10% | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1[b] | 1.5 | 2 | 1 | 1.5 | 2 | 1 | 1.5 | 2 |
| AICc | Q1 | 0.56 | 0.55 | 0.55 | 0.52 | 0.51 | 0.51 | 0.45 | 0.45 | 0.45 |
| | Q2 | 0.50 | 0.49 | 0.49 | 0.47 | 0.46 | 0.46 | 0.45 | 0.45 | 0.44 |
| | Q3 | 0.19 | 0.19 | 0.19 | 0.17 | 0.17 | 0.17 | 0.16 | 0.15 | 0.15 |
| | Q4 | 0.67 | 0.66 | 0.66 | 0.61 | 0.60 | 0.60 | 0.59 | 0.54 | 0.54 |
| | Q5 | 0.74 | 0.73 | 0.73 | 0.71 | 0.70 | 0.70 | 0.65 | 0.66 | 0.66 |
| LRT[c] | Q1 | 0.41 | 0.40 | 0.40 | 0.41 | 0.41 | 0.41 | 0.39 | 0.38 | 0.38 |
| | Q2 | 0.34 | 0.33 | 0.33 | 0.35 | 0.34 | 0.34 | 0.36 | 0.36 | 0.36 |
| | Q3 | 0.05 | 0.05 | 0.05 | 0.05 | 0.05 | 0.05 | 0.05 | 0.04 | 0.04 |
| | Q4 | 0.49 | 0.48 | 0.47 | 0.43 | 0.42 | 0.42 | 0.45 | 0.41 | 0.41 |
| | Q5 | 0.54 | 0.54 | 0.54 | 0.52 | 0.52 | 0.51 | 0.46 | 0.48 | 0.48 |

[a] Genome-wide significance level.

[b] 1, 1.5 and 2 are the amount subtracted from the LOD value at QTL position to estimate the LOD-$d$ support interval for the QTL.

[c] The critical value for the LRT at 5% significance level was obtained from the $\chi^2$ probability distribution with one degree of freedom.