

Bayesian QTL mapping using skewed Student-*t* distributions

Peter VON ROHR^{a,b}, Ina HOESCHELE^{a,*}

^a Departments of Dairy Science and Statistics,
Virginia Polytechnic Institute and State University,
Blacksburg, VA 24061-0315, USA

^b Institute of Animal Sciences, Animal Breeding,
Swiss Federal Institute of Technology (ETH), Zurich, Switzerland

(Received 23 April 2001; accepted 17 September 2001)

Abstract – In most QTL mapping studies, phenotypes are assumed to follow normal distributions. Deviations from this assumption may lead to detection of false positive QTL. To improve the robustness of Bayesian QTL mapping methods, the normal distribution for residuals is replaced with a skewed Student-*t* distribution. The latter distribution is able to account for both heavy tails and skewness, and both components are each controlled by a single parameter. The Bayesian QTL mapping method using a skewed Student-*t* distribution is evaluated with simulated data sets under five different scenarios of residual error distributions and QTL effects.

Bayesian QTL mapping / skewed Student-*t* distribution / Metropolis-Hastings sampling

1. INTRODUCTION

Most of the methods currently used in statistical mapping of quantitative trait loci (QTL) share the common assumption of normally distributed phenotypic observations. According to Coppieters *et al.* [2], these approaches are not suitable for analysis of phenotypes, which are known to violate the normality assumption. Deviations from normality are likely to affect the accuracy of QTL detection with conventional methods.

A nonparametric QTL interval mapping approach had been developed for experimental crosses (Kruglyak and Lander [8]) which was extended by Coppieters *et al.* [2] for half-sib pedigrees in outbred populations. Elsen and co-workers ([3, 7, 10]) presented alternative models for QTL detection in livestock populations. In a collection of papers these authors used heteroskedastic models

* Correspondence and reprints
E-mail: inah@vt.edu

to address the problem of non-normally distributed phenotypic observations. None of these methods can be applied to general and more complex pedigrees.

According to Fernandez and Steel [4], the existing toolbox for handling skewed and heavy-tailed data seems rather limited. These authors reviewed some of the existing approaches and concluded that they are all rather complicated to implement and lack flexibility and ease of interpretation.

Fernandez and Steel [4] have made an important contribution to the development of more flexible error distributions. They showed that by the method of inverse scaling of the probability density function on the left and on the right side of the mode, any continuous symmetric unimodal distribution can be skewed. This method requires a single scalar parameter, which completely determines the amount of skewness introduced into the distribution. This parameter must be estimated from the data. The procedure does not affect unimodality or tail behavior of the distribution. Simultaneously capturing heavy tails and skewness can be achieved by applying this method to a symmetric heavy-tailed distribution such as the Student- t distribution.

We believe that the approach developed by Fernandez and Steel [4] is one of the most promising methods to accommodate non-normal, continuous phenotypic observations with maximum flexibility. Fernandez and Steel [4] also demonstrated that this method is relatively easy to implement in a Bayesian framework. They designed a Gibbs sampler using data augmentation to obtain posterior inferences for a regression model with skewed Student- t distributed residuals.

The objective of this study was to incorporate the approach developed by Fernandez and Steel [4] into a Bayesian QTL mapping method, and to implement it with a Metropolis Hastings algorithm, instead of a Gibbs sampler with data augmentation, for better mixing of the Markov chain. In the following sections, we describe the method of inverse scaling, the QTL mapping model, a Markov chain Monte Carlo algorithm used to implement this method, and we show results from a simulation study. The simulated observations were generated from a model with one QTL flanked by two informative markers and a half-sib pedigree structure. Phenotypic error terms were assumed to follow four different distributions.

2. METHODS

2.1. Introducing skewness

In order to show how to introduce skewness into any symmetric and unimodal distribution, we closely followed the outline given by Fernandez and Steel [4]. Let us consider a univariate probability density function (pdf) $f(\cdot)$, which is unimodal and symmetric around 0. The pdf $f(\cdot)$ can be skewed by scaling the

density with inverse factors $\frac{1}{\gamma}$ and γ in the positive and negative orthant. This procedure will from now on be referred to as “inverse scaling of a pdf”, and it generates the following class of skewed distributions, indexed by γ :

$$p(e|\gamma) = \frac{2}{\gamma + \gamma^{-1}} \left\{ f\left(\frac{e}{\gamma}\right) I_{[0, \infty)}(e) + f(\gamma e) I_{(-\infty, 0)}(e) \right\} \quad (1)$$

where $\gamma \in \mathfrak{R}_+$ is a scalar, and $I_A(\cdot)$ stands for the indicator function over the set A .

For given values of γ and e , equation (1) specifies the probability density value for the skewed distribution associated with the specific value of γ . The term $f\left(\frac{e}{\gamma}\right)$ means that we have to evaluate the original symmetric pdf $f(\cdot)$ at value $\frac{e}{\gamma}$. Analogously, for $f(\gamma e)$, $f(\cdot)$ has to be evaluated at value γe . The indicator function can either take a value of 1, if the argument e to the function is within the set specified in the subscript of I , or a value of 0 otherwise. Factor $\frac{2}{\gamma + \gamma^{-1}}$ is a normalizing constant.

2.2. Properties of inverse scaling

The skewed pdf $p(e|\gamma)$ in (1) retains the mode at 0. From equation (1) it can be seen that the procedure of inverse scaling does not affect the location at which the maximum of the pdf occurs.

For $\gamma \neq 1$, the skewed pdf shown in equation (1) loses its symmetry. More formally this means that

$$p(e|\gamma \neq 1) \neq p(-e|\gamma \neq 1). \quad (2)$$

Inverting γ in equation (1) produces a mirror image around 0. Thus,

$$p(e|\gamma) = p\left(-e\left|\frac{1}{\gamma}\right.\right) \quad (3)$$

which in the case of $\gamma = 1$ leads to the property of symmetry.

The allocation of probability mass to each side of the mode is determined just by γ . This can also be seen from:

$$\frac{Pr(e \geq 0|\gamma)}{Pr(e < 0|\gamma)} = \gamma^2. \quad (4)$$

Fernandez and Steel [4] showed that the r -th order moment of (1) can be computed as:

$$E(e^r|\gamma) = M_r \frac{\gamma^{r+1} + \frac{-1^r}{\gamma^{r+1}}}{\gamma + \gamma^{-1}} \quad (5)$$

where

$$M_r = \int_0^{\infty} x^r 2f(x) dx.$$

The expression in (5) is finite, if and only if, the corresponding moment of the symmetric pdf $f(\cdot)$ exists.

Furthermore, Fernandez and Steel [4] gave a theorem which states that the existence of posterior moments for location and scale parameters in a linear model is completely unaffected by the added uncertainty of parameter γ . This means that these posterior moments exist, if and only if they also exist under symmetry where $\gamma = 1$.

2.3. Conditional distribution of phenotypes

In this section, we specify a Bayesian linear model for QTL mapping that accounts for skewness and heavy tails. Following the choice of Fernandez and Steel [4], we used the Student- t distribution as the symmetric pdf $f(\cdot)$. For a QTL mapping problem where phenotypes are assumed to be affected by a single QTL and a set of systematic factors, the model for trait values is as follows:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{T}_g\mathbf{v} + \mathbf{e} \quad (6)$$

where \mathbf{X} ($n \times r$) is design-covariate matrix, \mathbf{b} ($r \times 1$) is the vector of classification and regression effects, \mathbf{T}_g ($n \times q$) is the design matrix dependent on \mathbf{g} or the vector of QTL genotypes of all individuals, \mathbf{v} ($q \times 1$) is the vector of QTL effects, \mathbf{e} ($n \times 1$) is the vector of residuals, and n is the number of observations.

Here we assume that the QTL is bi-allelic, hence $q = 2$, $\mathbf{v} = [a, d]$, where a is half the difference between homozygotes and d is the dominance deviation. Row i of \mathbf{T}_g is $\mathbf{t}'_{i(g_i)} = [1, 0]$, $[0, 1]$, or $[-1, 0]$ if the individual i has QTL genotype $g_i = \text{QQ}$, Qq (or qQ) or qq , respectively.

Conditional on all unknown parameters and QTL genotypes, individual observations y_i are independent realizations from a distribution with probability density:

$$\begin{aligned} Pr(y_i | \mathbf{b}, \sigma_e^2, v, \gamma, a, d, g_i) &= \frac{2}{(\gamma + \gamma^{-1})} \frac{\Gamma\left(\frac{v+1}{2}\right)}{\Gamma\left(\frac{v}{2}\right) \sigma_e \sqrt{\pi v}} \\ &\times \left[1 + \frac{\left(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_{i(g_i)} \mathbf{v}\right)^2}{v \sigma_e^2} \right] \\ &\times \left\{ \frac{1}{\gamma^2} I_{[0, \infty)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_{i(g_i)} \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_{i(g_i)} \mathbf{v}) \right\}^{-\frac{v+1}{2}} \end{aligned} \quad (7)$$

where \mathbf{x}'_i is row i of matrix \mathbf{X} , and ν is the degrees-of-freedom parameter of the Student- t distribution.

The vector of unknowns in this problem is $(\mathbf{b}, \sigma_e^2, \nu, \gamma, a, d, p, \delta)$, where p denotes the QTL allele frequency and δ the genetic distance (in M assuming the Haldane mapping function) between one of the markers and the QTL. Note that model (6) depends on the vector of QTL genotypes, \mathbf{g} . Because of the simple pedigree structure, the likelihood of the phenotypes used in the Bayesian analysis was unconditional on the QTL genotypes, or

$$\begin{aligned} Pr(\mathbf{y}|\mathbf{b}, \sigma_e^2, \nu, \gamma, a, d, p, \delta) &= \prod_s^S \sum_{g_s} Pr(g_s|p) \\ &\times \prod_i^{n_s} \sum_{g_i} Pr(g_i|m_i, m_s, g_s; p, \delta) \\ &\times Pr(y_i|\mathbf{b}, \sigma_e^2, \nu, \gamma, a, d, g_i) \end{aligned} \quad (8)$$

where s denotes the father, S is the number of fathers, n_s is the number of offspring of the father s , g_s (g_i) is the QTL genotype of father s (offspring i), m_s (m_i) is the two-locus marker genotype of father s (offspring i) with phases assumed to be known, $Pr(g_s|p)$ is the Hardy-Weinberg frequency of genotype g_s which depends on QTL allele frequency p , and $Pr(g_i|m_i, m_s, g_s; p, \delta)$ depends on p (for the maternally inherited allele) and QTL position δ (for the paternally inherited allele).

The specific distribution of the error terms in model (6) introduces two additional parameters γ and ν into the problem.

2.4. Prior and posterior distributions

Different types of unknowns have independent prior distributions, or

$$\begin{aligned} Pr(\mathbf{b}, \sigma_e^2, \nu, \gamma, a, d, p, \delta) &= Pr(\mathbf{b}) \times Pr(\sigma_e^2) \times Pr(\nu) \times Pr(\gamma) \\ &\times Pr(a) \times Pr(d) \times Pr(p) \times Pr(\delta). \end{aligned} \quad (9)$$

For all unknowns, a uniform bounded prior was used. Such “uninformative” priors are appropriate in the absence of prior knowledge about the unknowns for specific traits, populations, and models as the one employed here. A list of prior distributions for all unknowns is given in Table I.

The joint posterior distribution of all unknowns was obtained (apart from a normalizing constant) by multiplying (9) with (8) using Table I.

Table I. Prior distributions for all unknowns used in the sampling scheme.

Unknown	Prior distribution	Hyper-parameter
b	Uniform $Pr(b) = \frac{1}{b_{\max} - b_{\min}}$	$b_{\min} = -5s_p$ $b_{\max} = 5s_p$
σ_e^2	Uniform $Pr(\sigma_e^2) = \frac{1}{\sigma_{e_{\max}}^2 - \sigma_{e_{\min}}^2}$	$\sigma_{e_{\min}}^2 > 0$ $\sigma_{e_{\max}}^2 < s_p^2$
ν	Uniform $Pr(\nu) = \frac{1}{\nu_{\max} - \nu_{\min}}$	$\nu_{\min} > 2$ $\nu_{\max} = s_p$
γ	Uniform $Pr(\gamma) = \frac{1}{\gamma_{\max} - \gamma_{\min}}$	$\gamma_{\min} > 0$ $\gamma_{\max} = s_p$
a	Uniform $Pr(a) = \frac{1}{a_{\max} - a_{\min}}$	$a_{\min} = -s_p$ $a_{\max} = s_p$
d	Uniform $Pr(d) = \frac{1}{d_{\max} - d_{\min}}$	$d_{\min} = -s_p$ $d_{\max} = s_p$
p	Uniform $Pr(p) = \frac{1}{p_{\max} - p_{\min}}$	$p_{\min} > 0$ $p_{\max} < 1$
δ	Uniform $Pr(\delta) = \frac{1}{\delta_{\max} - \delta_{\min}}$	$\delta_{\min} > 0$ $\delta_{\max} < 0.2$

s_p stands for the empirical phenotypic standard deviation of the observed data.

2.5. Metropolis Hastings (MH) sampling

The Metropolis Hastings algorithm was used to obtain samples from the joint posterior distribution of the parameters. With this algorithm and for a particular parameter, at each cycle t a candidate value y is proposed according to a proposal distribution $q(x, y)$, where x is the current sample value of the parameter. The candidate value is then accepted with probability $\alpha(x, y)$ where

$$\alpha(x, y) = \min\left(1, \frac{\pi(y)q(x, y)}{\pi(x)q(y, x)}\right) \quad (10)$$

and $\pi(\cdot)$ is the distribution one wants to sample from. Here, $\pi(\cdot)$ is the conditional distribution of an unknown parameter given the data and all

other unknowns. For a given unknown, the conditional distribution can be derived from the joint posterior distribution of all unknowns by retaining only those terms from the joint posterior which depend on the particular unknown. The conditional distributions for each unknown needed in (10) are given in Table II.

The proposal distributions $q(., .)$ were chosen to be uniform distributions centered at the current sample value with a small spread for all unknowns. The spread of the proposal distribution was determined by trial and error so that the overall acceptance rate of the samples was within the generally recommended range of [0.25, 0.4] (Chib and Greenberg [1]).

After a burn-in period of 2 000 cycles, an additional 100 000 cycles were generated. Posterior means of all unknowns were evaluated using all samples after the burn-in period. The length of the burn-in period was determined based on graphical inspection of the chains.

2.6. Simulation of data

Five scenarios of phenotypic distributions were considered. In the first scenario, the distribution of phenotypes was normal. This case represents a non-kurtosed symmetric error distribution. In the second scenario, we applied an inverse Box-Cox transformation, to this normal distribution, as described in MacLean *et al.* [9], to introduce skewness. A Student- t distribution, known to have heavy tails in the class of symmetric distributions, was used in the third scenario. In the fourth scenario, we employed a chi-square distribution, which is both kurtosed and skewed. Details about the distributions of the residuals used in the simulation are given in Table III. For these four scenarios, the phenotypes were influenced by a bi-allelic QTL with additive gene action and allele frequency of 0.5, which explained 12.5% of the phenotypic variation of the trait. The simulated pedigree had a half-sib structure with 40 sires each having 50 offspring. Because the focus of this study was on non-normal distributions of phenotypes rather than on how to deal with incomplete marker information, all fathers were heterozygous for the same pair of flanking markers and marker phases were assumed to be known. The distance between markers was 20 cM and the QTL was located at the midpoint of the marker interval.

Phenotypes under scenario five were simulated from the same χ^2 distribution as that used in scenario 4, but the effect of the QTL on the phenotype was set to zero. With this scenario we wanted to test whether the model would correctly predict that skewness in this case was not due to a putative QTL.

Vector \mathbf{b} contained the effects of one classification factor with three levels of -20 , 0 and 20 . Each data set was replicated 10 times.

Table II. Full conditional distributions for all unknowns using the priors in Table I.

(continued on the next page)

Unknown	Conditional distribution of the unknown given the data and all other unknowns
\mathbf{b}	$Pr(\mathbf{b} \sigma_e^2, \nu, \gamma, \mathbf{a}, \mathbf{d}, \mathbf{p}, \delta, \mathbf{y}) \propto \prod_{s=1}^S Pr(g_s p) \prod_{i=1}^{n_s} Pr(g_i m_i, m_s, g_s)$ $\times \left[1 + \frac{(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{\nu \sigma_e^2} \right]^{-\frac{\nu+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times \prod_{j=1}^k I_{[b_{\min}, b_{\max}]}(b_j)$
σ_e^2	$Pr(\sigma_e^2 \mathbf{b}, \nu, \gamma, \mathbf{a}, \mathbf{d}, \mathbf{p}, \delta, \mathbf{y}) \propto \sigma_e^{-n} \prod_{s=1}^S Pr(g_s p) \prod_{i=1}^{n_s} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{\nu \sigma_e^2} \right]^{-\frac{\nu+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times I_{[\sigma_{e\min}^2, \sigma_{e\max}^2]}(\sigma_e^2)$
ν	$Pr(\nu \mathbf{b}, \sigma_e^2, \gamma, \mathbf{a}, \mathbf{d}, \mathbf{p}, \delta, \mathbf{y}) \propto \left[\frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \right]^n (\nu)^{-\frac{n}{2}} \prod_{s=1}^S Pr(g_s p) \prod_{i=1}^{n_s} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{\nu \sigma_e^2} \right]^{-\frac{\nu+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times I_{[\nu_{\min}, \nu_{\max}]}(\nu)$
γ	$Pr(\gamma \mathbf{b}, \sigma_e^2, \nu, \mathbf{a}, \mathbf{d}, \mathbf{p}, \delta, \mathbf{y}) \propto \left[\frac{2}{\gamma + \gamma^{-1}} \right]^n \prod_{s=1}^S Pr(g_s p) \prod_{i=1}^{n_s} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{\nu \sigma_e^2} \right]^{-\frac{\nu+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(\mathbf{y}_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times I_{[\gamma_{\min}, \gamma_{\max}]}(\gamma)$

Table II. Continued.

Unknown	Conditional distribution of the unknown given the data and all other unknowns
a	$Pr(a \mathbf{b}, \sigma_e^2, v, \gamma, d, p, \delta, \mathbf{y}) \propto \prod_{s=1}^S \sum_{g_s} Pr(g_s p) \prod_{i=1}^{n_s} \sum_{g_i} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{v\sigma_e^2} \right]^{-\frac{v+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times \prod_{j=1}^k I_{[d_{\min}, d_{\max}]}(d_j)$
d	$Pr(d \mathbf{b}, \sigma_e^2, v, \gamma, a, p, \delta, \mathbf{y}) \propto \prod_{s=1}^S \sum_{g_s} Pr(g_s p) \prod_{i=1}^{n_s} \sum_{g_i} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{v\sigma_e^2} \right]^{-\frac{v+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times \prod_{j=1}^k I_{[d_{\min}, d_{\max}]}(d_j)$
p	$Pr(p \mathbf{b}, \sigma_e^2, v, \gamma, a, d, \delta, \mathbf{y}) \propto \prod_{s=1}^S \sum_{g_s} Pr(g_s p) \prod_{i=1}^{n_s} \sum_{g_i} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{v\sigma_e^2} \right]^{-\frac{v+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times \prod_{j=1}^k I_{[p_{\min}, p_{\max}]}(p_j)$
δ	$Pr(\delta \mathbf{b}, \sigma_e^2, v, \gamma, a, d, p, \mathbf{y}) \propto \prod_{s=1}^S \sum_{g_s} Pr(g_s p) \prod_{i=1}^{n_s} \sum_{g_i} Pr(g_i m_i, m_s, g_s; p, \delta)$ $\times \left[1 + \frac{(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v})^2}{v\sigma_e^2} \right]^{-\frac{v+1}{2}} \{ \gamma^{-2} I_{[0, \infty)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) + \gamma^2 I_{(-\infty, 0)}(y_i - \mathbf{x}'_i \mathbf{b} - \mathbf{t}'_i \mathbf{v}) \}$ $\times \prod_{j=1}^k I_{[\delta_{\min}, \delta_{\max}]}(\delta_j)$

Table III. Five different scenarios of simulating phenotypic distributions.

	Symmetric	Skewed	
Non-kurtosed	Normal	Skewed normal	
I	$[-20, 0, 20]'$	$[-20, 0, 20]'$	
Var $\langle e \rangle$	350	350	
a	10	10	
d	0	0	
p	0.5	0.5	
tp		0.1	
Kurtosed	Student-t	χ^2	χ^2 no QTL
I	$[-20, 0, 20]'$	$[-20, 0, 20]'$	$[-20, 0, 20]'$
Var $\langle e \rangle$	350	350	350
a	10	10	0
d	0	0	0
p	0.5	0.5	0
df	4	4	4

I stands for the vector of levels of the classification factor, a for half of the difference between homozygous QTL genotypes, d for the dominance deviation, p for the QTL allele frequency, tp for the transformation parameter described by McLean *et al.* [9], and df for the degrees of freedom of the Student- t and the χ^2 distribution used in the simulation.

3. RESULTS AND DISCUSSION

Tables IV–VIII summarize sample means, sample variances, Monte-Carlo standard errors (*MCSE*) and effective sample sizes (Geyer, [6]) for all unknowns. Sample means (sample variances) are averages across replicate data sets of the posterior means (variances) estimated from each Markov chain for individual parameters. *MCSE* is the square root of the variance of the average posterior mean estimate across replicates for a particular unknown. In Tables VII and VIII we also report averages across ten replicate data sets of posterior mean and variance for additive and dominance variance explained by the QTL.

Under the four scenarios which included a QTL in the simulation (Tabs. IV–VII), parameter estimates for the residual variance (Var $\langle e \rangle$), the QTL allele frequency (p), the QTL position (δ) and the three levels of the classification factor ($l_1 - l_3$) were close to their true values used in the simulation. The estimated QTL position δ was about 12 centimorgans from the left marker under all four scenarios that included a QTL, and significantly different from the true value for this parameter (10 cM) indicating a slight bias, which is not unusual for this type of QTL mapping analysis (see *e.g.* Zhang *et al.* [14]).

Table IV. Sample means^(a), sample variances^(b), Monte-Carlo standard errors ($MCSE$), and effective sample sizes^(c) ($EffSS$) for residual variance ($\text{Var}(e)$), degrees of freedom parameter (ν), skewness parameter (γ), half of the difference between homozygotes (a), dominance deviation (d), QTL allele frequency (p), QTL position (δ), and three levels of the classification factor (l_1 , l_2 and l_3) under the normal scenario.

	True value	Sample mean	Sample variance	$MCSE$	$EffSS$
Scenario	normal				
$\text{Var}(e)$	350	315.4	504.6	1.041	1 734
ν	∞	17.16	4.851	0.0291	6 993
γ	1	1.006	0.0020	0.0014	1 045
a	10	7.680	4.850	0.1111	1 328
d	0	6.998	38.38	0.4479	336
p	0.5	0.5159	0.0089	0.0040	1 079
δ	0.1	0.1196	0.0002	0.0001	13 780
l_1	-20	-20.71	8.257	0.2092	302
l_2	0	-0.7200	8.306	0.2093	297
l_3	20	19.07	8.218	0.2092	293

^(a) Average across replicate data sets, posterior mean estimate.

^(b) Average across replicate data sets, posterior variance estimate.

^(c) As calculated in Geyer [6].

Table V. Sample means^(a), sample variances^(b), Monte-Carlo standard errors ($MCSE$), and effective sample sizes^(c) ($EffSS$) for residual variance ($\text{Var}(e)$), degrees of freedom parameter (ν), skewness parameter (γ), half of the difference between homozygotes (a), dominance deviation (d), QTL allele frequency (p), QTL position (δ), and three levels of the classification factor (l_1 , l_2 and l_3) under the skewed-normal scenario.

	True value	Sample mean	Sample variance	$MCSE$	$EffSS$
Scenario	skewed-normal				
$\text{Var}(e)$	350	349.5	432.0	0.9466	1 023
ν	∞	16.95	5.118	0.0267	8 240
γ		1.430	0.0052	0.0030	664
a	10	7.364	12.28	0.4112	520
d	0	6.290	42.91	0.6466	280
p	0.5	0.4830	0.0085	0.0070	1 072
δ	0.1	0.1212	0.0002	0.0001	12 708
l_1	-20	-18.49	10.76	0.3207	284
l_2	0	1.395	10.70	0.3190	285
l_3	20	21.13	10.68	0.3187	291

^(a) Average across replicate data sets, posterior mean estimate.

^(b) Average across replicate data sets, posterior variance estimate.

^(c) As calculated in Geyer [6].

Table VI. Sample means^(a), sample variances^(b), Monte-Carlo standard errors (*MCSE*), and effective sample sizes^(c) (*EffSS*) for residual variance ($\text{Var}(e)$), degrees of freedom parameter (ν), skewness parameter (γ), half of the difference between homozygotes (a), dominance deviation (d), QTL allele frequency (p), QTL position (δ), and three levels of the classification factor (l_1 , l_2 and l_3) under the Student- t scenario.

	True value	Sample mean	Sample variance	<i>MCSE</i>	<i>EffSS</i>
Scenario			Student-t		
$\text{Var}(e)$	350	321.2	557.0	0.1882	16 899
ν	4	4.340	0.2493	0.0068	5 527
γ	1	1.021	0.0014	0.0009	1 983
a	10	9.587	1.519	0.0443	2 249
d	0	1.911	8.250	0.1381	860
p	0.5	0.4991	0.0063	0.0027	1 398
δ	0.1	0.1222	0.0002	0.0001	14 208
l_1	-20	-19.84	2.283	0.0629	789
l_2	0	0.7000	2.3286	0.0613	818
l_3	20	20.26	2.287	0.0615	803

(a) Average across replicate data sets, posterior mean estimate.

(b) Average across replicate data sets, posterior variance estimate.

(c) As calculated in Geyer [6].

Under the scenarios with the Student- t and the χ^2 distribution with a QTL, the estimates for a and d were close to the true values used in the simulation, and the sample variances and *MCSE* were lower than under the other scenarios. For the normal and skewed normal distributions, a and d were estimated less accurately, and sample variances and *MCSE* were higher (to some extent, this also applies to parameter p).

The estimates for parameters a and d under the scenario with the χ^2 distribution without a QTL (Tab. VIII) deviated from their true values of zero. Posterior variances and *MCSE* of these parameters were very high, and effective sample sizes were extremely small, with similar results for the other location parameters (the three levels of the classification factor), indicating poor identifiability of these parameters.

To see whether our method can effectively discriminate between a non-normal phenotypic distribution with a QTL (χ^2) and a non-normal distribution without a QTL (χ^2 no QTL), we first estimated the marginal posterior densities of the additive ($2p(1-p)[a+d(p-q)]^2$) and dominance ($4p^2(1-p)^2d^2$) variances of the QTL shown as histograms for one replicate data set under the χ^2 scenario with QTL in Figure 1 and under the χ^2 scenario without QTL in Figure 2. The histograms show a very high frequency for an additive QTL

Table VII. Sample means^(a), sample variances^(b), Monte-Carlo standard errors (*MCSE*), and effective sample sizes^(c) (*EffSS*) for residual variance ($\text{Var}(e)$), degrees of freedom parameter (ν), skewness parameter (γ), half of the difference between homozygotes (a), dominance deviation (d), QTL allele frequency (p), QTL additive variance (σ_a^2), QTL dominance variance (σ_d^2), QTL position (δ), and classification factor (l_1, l_2 and l_3) under the χ^2 scenario.

	True value	Sample mean	Sample variance	<i>MCSE</i>	<i>EffSS</i>
Scenario			χ^2		
$\text{Var}(e)$	350	331.4	296.4	0.1875	12 295
ν		11.80	7.060	0.0393	4 997
γ		3.179	0.1390	0.0220	322
a	10	9.377	0.3367	0.0152	2 633
d	0	0.7039	0.5610	0.0205	2 404
p	0.5	0.4963	0.0017	0.0009	2 931
σ_a^2	50	43.62	30.25	0.0292	3 597
σ_d^2	0	0.3001	0.4422	0.0016	3 196
δ	0.1	0.1139	0.0001	0.0001	14 173
l_1	-20	-20.47	0.6536	0.0200	2 490
l_2	0	-0.572	0.6767	0.0206	2 576
l_3	20	19.17	0.6111	0.0179	2 818

^(a) Average across replicate data sets, posterior mean estimate.

^(b) Average across replicate data sets, posterior variance estimate.

^(c) As calculated in Geyer [6].

variance close to 0 under the scenario without a QTL, whereas under the scenario with a QTL, 0 was not within the displayed range. The frequency for the dominance QTL variance was highest around the true value of 0 under the scenario with a QTL. Under the scenario without QTL, the maximum frequency occurred at a higher variance value, and the range of the QTL dominance variance was larger.

From the marginal posterior distributions, we also estimated the boundaries of 95% Highest Posterior Density (HPD) regions as described by Tanner [12]. Average boundaries across ten replicate data sets were 18.46 and 67.30 for the QTL additive variance, and 0.089 and 8.287 for the QTL dominance variance under the χ^2 scenario with a QTL. Under the χ^2 scenario without a QTL the boundaries were 0.000 and 262.4 for the QTL additive and 0.000 and 44.07 for the QTL dominance variance. The boundaries of the HPD regions included the value of zero for the QTL additive variance in five out of ten replicate data sets under the scenario without a QTL, and for the five other replicates, the lower boundary of the HPD region was very close to zero (average lower boundary

Table VIII. Sample means^(a), sample variances^(b), Monte-Carlo standard errors (*MCSE*), and effective sample sizes^(c) (*EffSS*) for residual variance ($\text{Var}(e)$), degrees of freedom parameter (ν), skewness parameter (γ), half of the difference between homozygotes (a), dominance deviation (d), QTL allele frequency (p), QTL additive variance (σ_a^2), QTL dominance variance (σ_d^2), QTL position (δ), and classification factor (l_1 , l_2 and l_3) under the χ^2 scenario no QTL.

	True value	Sample mean	Sample variance	<i>MCSE</i>	<i>EffSS</i>
Scenario			χ^2 no QTL		
$\text{Var}(e)$	350	327.0	349.1	1.258	1 374
ν		11.05	6.489	0.1042	1 796
γ		5.215	1.351	0.1510	67
a	0	2.228	40.62	1.645	17
d	0	6.325	22.66	1.120	31
p		0.4896	0.0406	0.0438	26
σ_a^2	0	20.45	1595	245.2	38.43
σ_d^2	0	7.081	29.61	1.840	45.90
δ		0.1229	0.0003	0.0001	16 152
l_1	-20	-15.18	21.71	1.131	35
l_2	0	4.821	22.04	1.138	36
l_3	20	24.71	21.30	1.119	36

^(a) Average across replicate data sets, posterior mean estimate.

^(b) Average across replicate data sets, posterior variance estimate.

^(c) As calculated in Geyer [6].

was 5.11 for these five replicates). For the scenario with a QTL, the value of zero was included in the HPD region for the additive QTL variance only in one out of ten replicates. The true value for the QTL additive variance of 50 was within the HPD region for every replicate under the scenario with a QTL. The HPD region for the QTL dominance variance was much wider under the scenario without a QTL compared to the scenario with a QTL. The HPD regions for the dominance variance included the true value of zero in seven (eight) out of ten replicates for the χ^2 with QTL (without QTL) scenario.

All data sets representing the χ^2 distribution scenarios were analyzed with a model that assumes normal phenotypes. Under both scenarios (with and without a QTL), residual, additive QTL and dominance QTL variance estimates were much closer to the true value when the analysis was performed with the skewed Student-*t* model rather than with the normal model. Assuming normal phenotypes under the two χ^2 scenarios caused the residual variance to be underestimated, while additive and dominance QTL variance were both overestimated considerably (Tab. IX). The HPD regions for the QTL additive

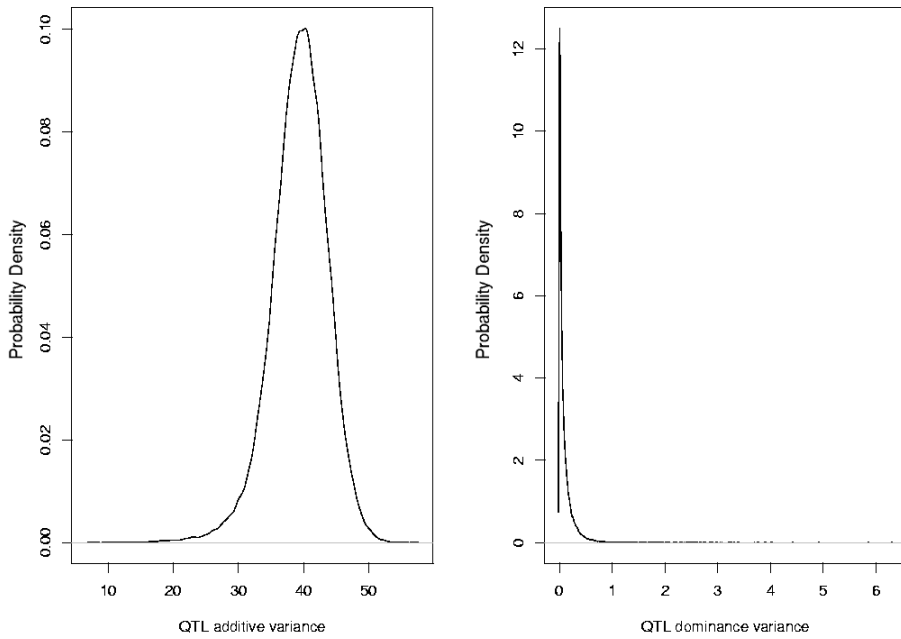


Figure 1. Marginal posterior densities of QTL additive and QTL dominance variance under the χ^2 scenario with a QTL.

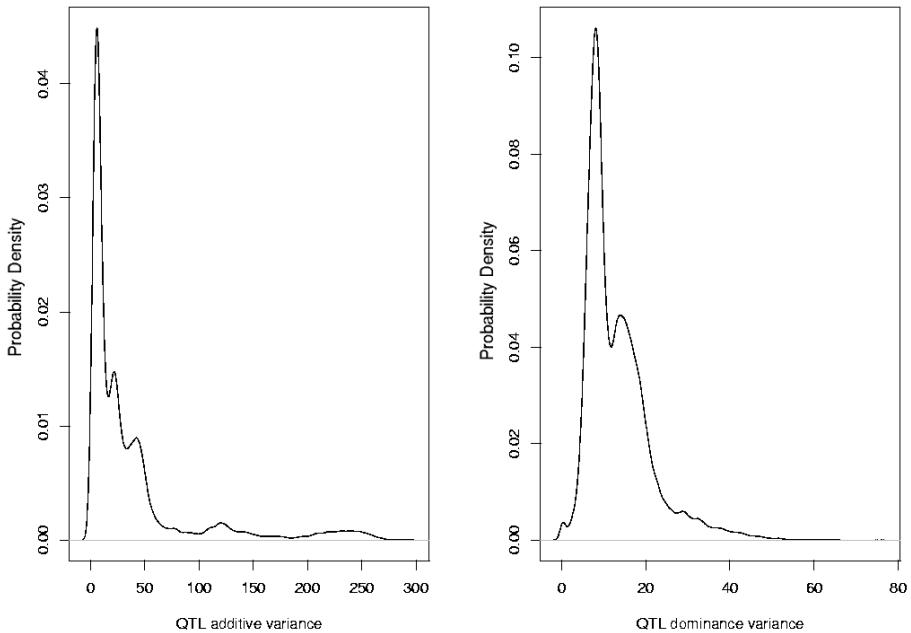


Figure 2. Marginal posterior densities of QTL additive and QTL dominance variance under the χ^2 scenario without a QTL.

Table IX. Sample means^(a), sample variances^(b), Monte-Carlo standard errors (*MCSE*), and effective sample sizes^(c) (*EffSS*) for residual variance ($\text{Var}(e)$), QTL additive variance (σ_a^2), and QTL dominance variance (σ_d^2) under both χ^2 scenarios (with and without a QTL) analyzed with a normal penetrance function.

	True value	Sample mean	Sample variance	<i>MCSE</i>	<i>EffSS</i>
Scenario			χ^2		
$\text{Var}(e)$	350	208.9	129.7	0.0151	9 097
σ_a^2	50	148.4	162.6	0.0430	4 458
σ_d^2	0	28.22	26.76	0.0038	5 998
Scenario			χ^2 no QTL		
$\text{Var}(e)$	350	179.4	59.98	0.0040	15 220
σ_a^2	0	159.6	37.19	0.0030	10 156
σ_d^2	0	26.49	10.96	0.0012	11 266

(a) Average across replicate data sets, posterior mean estimate.

(b) Average across replicate data sets, posterior variance estimate.

(c) As calculated in Geyer, [6].

variance contained a true value of 50 in none of the replicates for the scenario with the QTL and the value of 0 in five out of ten replicates for the scenarios with and without a QTL. The true value of 0 for the QTL dominance variance was outside of the HPD regions in all replicates for both scenarios with and without a QTL when analyzing the data with a normal penetrance function.

These results indicate that we would detect the absence of a QTL 50% of the time, when we only consider inclusion of the value of zero in the HPD region for the QTL additive variance. However, in the absence of a QTL, the lower boundary of the HPD region always either included zero or was close to zero, and the HPD region was very wide, indicating little information and support for a QTL. Replacement of the normal by the skewed Student-*t* penetrance function clearly improved the accuracy of parameter estimation.

A value of the skewness parameter (γ) close to 1 indicates a symmetric distribution. This was the case for the normal and the Student-*t* distribution. Estimates for γ were 1.006 under the normal and 1.021 under the Student-*t* scenario. Under the three scenarios with skewed error distributions, estimates for γ ranged between 1.430 and 5.215, and thus indicated the presence of skewness in the distribution of residual phenotypes.

Parameter ν represents the degrees of freedom under a Student-*t* distribution with symmetry ($\gamma = 1$). In our simulations, we used four degrees of freedom under the Student-*t* scenario. With a value of 4.340 the estimate of ν was close to the true value. Under a skewed Student-*t* distribution with $\gamma \neq 1$,

parameter ν is a measure of the tail behavior. The smaller the ν , the heavier were the tails of the distribution. Based on the estimates of ν , the scenarios used in the simulation can be categorized as heavy tailed such as the Student- t or not heavy-tailed such as the normal and the skewed-normal showing larger estimates of ν . Under the two scenarios with the χ^2 residuals, estimates of ν were in-between the estimates from the Student- t and the normal distribution.

In a previous study (von Rohr and Hoeschele [13]), we reported that estimates of ν are somewhat dependent on the prior distribution for ν . In this study we chose a bounded uniform prior distribution for ν . In theory, the value of ν tends to infinity for the normal distribution. Hence, although the range of the bounded uniform prior distribution does not cover normal distributions, high estimates of ν (near the upper bound) are obtained when residual phenotypes are normally distributed, and thus indicate little deviation from normality.

Posterior correlations between parameters were estimated from the sample values of the Markov chains and are listed in Table X. The strongest correlations were obtained between the QTL parameters defining the variance explained by the QTL (a, d, p) and between all phenotypic mean parameters (a, d, l_1, l_2, l_3). A comparison of the correlations between scenarios showed that they tended to be lower under the symmetric distributions (Student- t and normal) than under skewed error distributions.

4. CONCLUSIONS

A robust Bayesian QTL mapping method was implemented, which allows for non-normal, continuous distributions of phenotypes within QTL genotypes, via skewed Student- t distributions of residual phenotypes in the analysis. The skewed Student- t distribution was obtained by the method of inverse scaling, and this approach can handle distributions where skewness or heavy tails or both are present. Overall, this study confirms the good results reported by Fernandez and Steel [4], who showed that this method can handle even more extreme cases such as the stable distribution. Parameters were estimated with good accuracy under a range of distributions, except for for the normal distribution where additive QTL effects were underestimated and dominance effects overestimated. Hence, if ν and γ parameters indicate no deviation from the normal distribution (as was the case under the normal scenario here), one should reanalyze the data with the normal penetrance function to obtain more accurate parameter estimates (as we confirmed, results not shown). When there is deviation from normality, parameters should be estimated more accurately with the skewed Student- t than with the normal penetrance function, as we demonstrated for the χ^2 -distribution with QTL. There did not appear to be much of a difference between analyses using normal or skewed Student- t penetrance functions, when applied to a skewed and kurtosed distribution without a QTL, in

the indication of QTL absence or little support for a QTL. However, parameter estimation was much improved by use of the skewed Student- t penetrance function.

Fernandez and Steel [4] and Strandén and Gianola [11], among others, used a Gibbs sampler with data augmentation to sample from the joint posterior distribution for problems involving Student- t distributions. Data augmentation was motivated by the representation of the Student- t distribution as a scale mixture of normals. Data augmentation facilitates sampling by producing standard conditional distributions which are convenient to sample from. Data augmentation comes at the expense of an additional mixing parameter λ_i for each observation i . We implemented a Metropolis-Hastings sampler, which resulted in a simple sampling scheme and has the advantages of avoiding data augmentation and controlling autocorrelations among successive samples to some extent via choice of proposal distributions. The performance of the method of inverse scaling, *i.e.* the replacement of the normal by the skewed Student- t penetrance function, in the simple QTL model considered here indicates that this approach should also be useful for more complex QTL models including multiple QTLs and complex pedigrees. Applying this approach to complex pedigrees would include fitting a residual polygenic effect. Strandén and Gianola [11] proposed to use a symmetric Student- t distribution for polygenic effects. Their results did not indicate that Student- t distributed polygenic effects would be beneficial to the analysis.

ACKNOWLEDGEMENTS

This work was supported by the *National Science Foundation* grant DBI-9723022 to Ina Hoeschele, a fellowship from the Swiss National Science Foundation to Peter von Rohr, and the *National Center for Supercomputing Applications* under grant number MCB990003N and utilized the computer system *SGI Origin2000* at the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign.

Comments from two anonymous reviewers were greatly appreciated.

REFERENCES

- [1] Chib S., Greenberg E., Understanding the Metropolis-Hastings algorithm, *Amer. Stat.* 49 (1995) 327–335.
- [2] Coppieters W., Kvasz A., Farnir F., Arranz J.J., Grisart B., Mackinnon M., Georges M., A Rank-based nonparametric method for mapping quantitative trait loci in outbred half-sib pedigrees: Application to milk production in a granddaughter design, *Genetics* 149 (1998) 1547–1555.

- [3] Elsen J.-M., Mangin B., Goffinet B., Boichard D., Le Roy P., Alternative models for QTL detection in livestock. I. General introduction, *Genet. Sel. Evol.* 31 (1999) 213–224.
- [4] Fernandez C., Steel MFJ, On Bayesian modeling of fat tails and skewness, *J. Am. Statist. Assoc.* 93 (1998) 359–371.
- [5] Geweke J., Bayesian treatment of the independent Student-*t* linear model, *J. Appl. Econometrics* 8 (1993) S19–S40.
- [6] Geyer C.J., Practical Markov Chain Monte Carlo, *Stat. Sci.* 7 (1992) 473–511.
- [7] Goffinet B., Le Roy P., Boichard D., Elsen J.-M., Mangin B., Alternative models for QTL detection in livestock. III. Heteroskedastic model and models corresponding to several distributions of the QTL effect, *Genet. Sel. Evol.* 31 (1999) 341–350.
- [8] Kruglyak L., Lander E.S., A Nonparametric approach for mapping quantitative trait loci, *Genetics* 139 (1995) 1421–1428.
- [9] MacLean C.J., Morton N.E., Elston R.C., Yee S., Skewness in commingled distributions, *Biometrics* 32 (1976) 695–699.
- [10] Mangin B., Goffinet B., Le Roy P., Boichard D., Elsen J.-M., Alternative models for QTL detection in livestock. II. Likelihood approximations and sire marker genotype estimation, *Genet. Sel. Evol.* 31 (1999) 225–237.
- [11] Strandén I., Gianola D., Mixed effects linear models with *t*-distributions for quantitative genetic analysis: A Bayesian approach, *Genet. Sel. Evol.* 31 (1999) 25–42.
- [12] Tanner MA., *Tools for Statistical Inference*, 2nd edn., in: Springer Series in Statistics, New York, 1993.
- [13] von Rohr P., Hoeschele I., Robust Bayesian analysis using skewed Student-*t* distributions, in: 50th Annual Meeting of EAAP, August 22–26, Zurich, paper G3.7.
- [14] Zhang Q., Boichard D., Hoeschele I., Ernst C., Eggen A., Murkve B., Pfister-Genskow M., Witte L.A., Grignola F.E., Uimari P., Thaller G., Bishop M.D., Mapping quantitative trait loci for milk production and health of dairy cattle in a large outbred pedigree, *Genetics* 149 (1998) 1959–1973.

Table X. Continued.

	ν	γ	a	d	p	δ	l_1	l_2	l_3
	χ^2								
Var $\langle e \rangle$	-0.397	0.146	-0.089	0.107	-0.002	0.025	0.285	0.276	0.300
ν		0.198	-0.064	0.075	0.027	0.022	0.030	0.029	0.048
γ			0.037	0.023	-0.059	-0.014	-0.099	-0.098	-0.079
a				0.466	0.132	0.058	0.221	0.234	0.217
d					0.094	0.032	0.419	0.463	0.425
p						-0.006	0.187	0.190	0.158
δ							-0.006	-0.002	0.001
l_1								0.648	0.632
l_2									0.647
	χ^2 no QTL								
Var $\langle e \rangle$	-0.445	0.142	-0.017	0.261	-0.010	-0.002	-0.155	-0.151	-0.152
ν		0.061	-0.058	0.029	-0.044	-0.006	0.018	0.023	0.017
γ			0.101	0.238	0.093	0.001	-0.301	-0.287	-0.290
a				0.524	0.829	0.010	0.594	0.594	0.593
d					0.376	0.018	0.937	0.937	0.938
p						-0.003	-0.075	-0.077	-0.075
δ							-0.007	-0.008	-0.007
l_1								0.984	0.983
l_2									0.985

^(a) Posterior correlations estimated as sample correlations from *MCMC* output and averaged across ten replicate data sets.