

A note on QTL detecting for censored traits

Yixin FANG*

Department of Statistics, Columbia University, New York, NY 10027, USA

(Received 9 August 2005; accepted 6 October 2005)

Abstract – Most existing statistical methods for mapping quantitative trait loci (QTL) assume that the phenotype follows a normal distribution and that it is fully observed. However, some phenotypes have skewed distributions and may be censored. This note proposes a simple and efficient approach to QTL detecting for censored traits with the Cox PH model without estimating the baseline hazard function which is “nuisance”.

QTL detecting / censored trait / interval mapping

1. INTRODUCTION

The standard approach for mapping the quantitative trait loci (QTL) contributing to variation in a quantitative trait makes use of the assumption that the phenotype is normally distributed and fully observed [7,8,11,18]. In recent years, many authors have proposed nonparametric or semiparametric methods to solve the problem of model misspecification caused by the assumption of normality [6,10,19]. In addition, assumptions of standard approach are likely to be false when the phenotype pertains to the survival time or failure time and the failure time is often subject to censoring. The incompleteness of the trait values presents a major challenge in the application of the interval-mapping approach [11]. Symons *et al.* [17] computed LOD scores under the Cox proportional hazards (PH) model using a variant of the EM algorithm using Monte Carlo simulation to make the computations tractable as described by Lipsitz and Ibrahim [13]. The EM algorithm incorporates all possible values of missing covariates according to the appropriate probability distributions. This method is computationally intensive and it needs estimating a “nuisance” parameter, that is, baseline hazard function $\lambda_0(\cdot)$ of the Cox PH model. Broman [2] considered a cure-rate model in which the mice

* Corresponding author: yf2113@columbia.edu

that are alive at the end of the study are regarded as cured and in which the survival times among the deaths follow a log-normal distribution. This is a special model which can only deal with the situations in which the potential censoring times are equal among all study subjects. Diao *et al.* [4] formulated the effects of QTL on the failure time through a parametric PH model with a Weibull baseline hazard function $\lambda_0(t) = \gamma_1\gamma_2 t^{\gamma_2-1}$, $\gamma_1 > 0$, $\gamma_2 > 0$. Since it is a parametric model, it still has the problem of mis-specification. In a recent study, Moreno *et al.* [15] proposed two new QTL detection approaches which allow the consideration of censored data. One is similar to Diao *et al.* [4] based on Weibull distribution and the other one is based on Cox proportional hazards model.

Since the primary reason for using the Cox proportional hazards model and his partial likelihood technique is avoidance of the “nuisance” baseline hazard function $\lambda_0(\cdot)$, in this article, we provide a simple interval-mapping method to censor traits without estimating $\lambda_0(\cdot)$. In brief, we formulate the effects of QTL on the failure time through the PH model and treat unobserved genotypes of QTL as missing covariates. Then we develop a procedure based on partial likelihood for detecting QTL and show how to assess genome-wide statistical significance. In comparison to Symons *et al.* [17], Broman [2] and Diao *et al.* [4], our approach has some advantages. First, we used the Cox semi-parametric PH model which is most popular for survival analysis. Second, we avoided estimating a baseline hazard function which is very complicated to be estimated. Third, the test statistic we used was actually the well-known log rank statistic and it is the locally most powerful test. Furthermore, because there was no iteration in calculating the test statistic, the method proposed in the following was computationally efficient.

2. MAIN RESULTS

We consider populations derived from a cross between two parental inbred lines P_1 and P_2 . There are two kinds of basic populations, F_2 and backcross. In this note we are only concerned with F_2 . Consider n progenies from an F_2 population. Let T_i denote the quantitative trait from the i th subject, which pertains to a failure time that can potentially be censored and thus incompletely observed. Let C_i be the censoring time for the i th subject. The observation consists of $\tilde{T}_i = \min(T_i, C_i)$ and $\delta_i = I(T_i \leq C_i)$, where $I(A)$ is the indicator function for event A . The failure time T_i is fully observed only when it is uncensored, *i.e.*, $\delta_i = 1$.

Suppose that we have data on a set of genetic markers with a known genetic map. Let M_i denote the multipoint marker genotype data for the

i th subject. We consider a putative QTL locus d in the genome and define $G_i = -1, 0$ or 1 according to whether the i th subject has genotype qq, Qq or QQ , respectively, at the QTL. We specify a proportional hazards model for the effects of the QTL genotypes on the failure time such that, conditional on the QTL genotype G_i , the hazard function of T_i takes the form

$$\lambda(t|G_i) = \lambda_0(t) \exp\{\beta_1 G_i + \beta_2(1 - |G_i|)\}, i = 1, \dots, n, \tag{1}$$

where β_1 and β_2 pertain to the additive and dominant effects of QTL and $\lambda_0(t)$ is an unknown baseline hazard function. Diao *et al.* [4] considered a Weibull hazard function $\lambda_0(t) = \gamma_1 \gamma_2 t^{\gamma_2 - 1}, \gamma_1 > 0, \gamma_2 > 0$. In this article, we add no condition on the form of $\lambda_0(t)$.

Because G_i 's are missing covariates in model (1), we consider conditional hazard function given \mathbf{M}_i , that is

$$\lambda(t|\mathbf{M}_i) = \lambda_0(t) E \left[e^{\beta_1 G_i + \beta_2(1 - |G_i|)} | \mathbf{M}_i, T_i \geq t \right], \tag{2}$$

which is also a multiplicative hazards model. Denote the conditional expectation in (2) by $a_i(t, \boldsymbol{\beta})$, where $\boldsymbol{\beta} = (\beta_1, \beta_2)'$, then we have

$$a_i(t, \boldsymbol{\beta}) = \frac{E \left[e^{\beta_1 G_i + \beta_2(1 - |G_i|)} \exp \left\{ -\Lambda_0(t) e^{\beta_1 G_i + \beta_2(1 - |G_i|)} \right\} | \mathbf{M}_i \right]}{E \left[\exp \left\{ -\Lambda_0(t) e^{\beta_1 G_i + \beta_2(1 - |G_i|)} \right\} | \mathbf{M}_i \right]}, \tag{3}$$

where $\Lambda_0(t) = \int_0^t \lambda_0(s) ds$ is the cumulative baseline hazard function.

Suppose now that $t_1 < \dots < t_k$ are the ordered distinct failures in the sample and $R(t_j)$ and $D(t_j)$ denote the risk set just prior to t_j and the set of subjects failing at t_j , respectively, $j = 1, \dots, k$. If we use an approximation to accommodate tied failure times [1], like Prentice [16], the partial likelihood function is given as

$$L(\boldsymbol{\beta}) = \prod_{j=1}^k \frac{\prod_{h \in D(t_j)} a_h(t_j, \boldsymbol{\beta})}{\left[\sum_{h \in R(t_j)} a_h(t_j, \boldsymbol{\beta}) \right]^{m_j}}, \tag{4}$$

where m_j is the number of failures at t_j and n_j is the size of the risk set $R(t_j)$, $j = 1, \dots, k$.

The score function of (4) is

$$\mathbf{S}(\boldsymbol{\beta}) = \frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{j=1}^k \left\{ \sum_{h \in D(t_j)} \frac{\mathbf{b}_h(t_j, \boldsymbol{\beta})}{a_h(t_j, \boldsymbol{\beta})} - \frac{m_j \sum_{h \in R(t_j)} \mathbf{b}_h(t_j, \boldsymbol{\beta})}{\sum_{h \in R(t_j)} a_h(t_j, \boldsymbol{\beta})} \right\}, \tag{5}$$

where $\mathbf{b}_i(t, \boldsymbol{\beta}) = \partial a_i(t, \boldsymbol{\beta}) / \partial \boldsymbol{\beta}, i = 1, \dots, n$.

In order to test the null hypothesis of no QTL effect, *i.e.*, $H_0 : \boldsymbol{\beta} = \mathbf{0}$, we used score test procedure. To do so, let $\mathbf{x}_i = (x_{1i}, x_{2i})'$, $i = 1, \dots, n$, where $x_{1i} = E[G_i | \mathbf{M}_i]$ and $x_{2i} = E[(1 - |G_i|) | \mathbf{M}_i]$. These values, x_{1i} 's and x_{2i} 's, can be found, for example, in [14]. It is easy to verify that $a_i(t, \mathbf{0}) = 1$ and $\mathbf{b}_i(t, \mathbf{0}) = \mathbf{x}_i$, for any $i = 1, \dots, n$, then the score statistic (5) at $\boldsymbol{\beta} = \mathbf{0}$ can be written as

$$\mathbf{s} = \sum_{j=1}^k \left(\sum_{h \in D(t_j)} \mathbf{x}_h - m_j n_j^{-1} \sum_{h \in R(t_j)} \mathbf{x}_h \right). \tag{6}$$

We were not excited to find that the above statistic does not depend on the nuisance parameter $\lambda_0(\cdot)$. Since the primary reason for using the partial likelihood technique is avoidance of $\lambda_0(\cdot)$, the use of statistic (6) will lead to a simple and efficient mapping approach.

By some arguments based on the counting process (App. A), we can get an estimate of variance of score statistics (6),

$$\mathbf{v} = \sum_{j=1}^k m_j n_j^{-1} \left[\sum_{h \in R(t_j)} \mathbf{x}_h^{\otimes 2} - n_j^{-1} \left(\sum_{h \in R(t_j)} \mathbf{x}_h \right)^{\otimes 2} \right], \tag{7}$$

where $\mathbf{x}^{\otimes 2} = \mathbf{x}\mathbf{x}'$ for \mathbf{x} a vector.

In fact, similar to Prentice [16], a finite population variance argument applied to $\sum_{h \in D(t_j)} \mathbf{x}_h$, given $R(t_j)$ for each $j = 1, \dots, k$ leads to

$$\tilde{\mathbf{v}} = \sum_{j=1}^k \frac{n_j - m_j}{n_j - 1} \cdot \frac{m_j}{n_j} \left[\sum_{h \in R(t_j)} \mathbf{x}_h^{\otimes 2} - n_j^{-1} \left(\sum_{h \in R(t_j)} \mathbf{x}_h \right)^{\otimes 2} \right], \tag{8}$$

where factors $\frac{n_j - m_j}{n_j - 1}$ due to the tied data approximation used in (4).

Under $H_0 : \boldsymbol{\beta} = \mathbf{0}$, the statistic

$$\mathbf{w} = \mathbf{s}' \tilde{\mathbf{v}}^{-1} \mathbf{s}, \tag{9}$$

will have an asymptotic χ^2_2 distribution. Note that the score \mathbf{s} and variance \mathbf{v} or $\tilde{\mathbf{v}}$ all depend on the locus d of QTL through the dependence of \mathbf{x}_i 's on d . In the sequel, we include d in the expressions to emphasize their dependence on d ,

$$\mathbf{w}(d) = \mathbf{s}'(d) \tilde{\mathbf{v}}^{-1}(d) \mathbf{s}(d). \tag{10}$$

Thus the test statistic curve $\{\mathbf{w}(d), d \geq 0\}$ for each chromosome can be drawn as in the case of standard interval mapping. For each chromosome, the position with the largest value of the curve is declared to be the QTL location provided that the value exceeds a certain threshold level. In the next section, we will show how to determine the threshold level.

3. THRESHOLD VALUES

When searching the entire chromosome or whole genome for QTL, one should select a threshold level such that the probability that the test statistic exceeds this level anywhere in the genome equals the desired false-positive rate. In Appendix B we show that in a dense-map case, the process $\{\tilde{v}^{-1/2}(d)s(d), d \geq 0\}$ is asymptotically equivalent to a two dimensional Ornstein-Uhlenbeck process under a null hypothesis. Thus we can get analytical approximations of thresholds which are analogous to those of Lander and Botstein [11], Dupuis and Siegmund [5], etc.

Since the analytical results are based on a number of assumptions that are not likely to be met in practice, we can simply use a permutation test, as described by Churchill and Doerge [3], to obtain an empirical threshold value. In addition, this section concludes with a resampling procedure similar to Diao *et al.* [4] and Zou *et al.* [20] by which we approximate the null distribution of $\sup_d w(d)$ and then get the threshold value of our interval mapping method. First, we generate $Z_i, i = 1, \dots, n$, which are *i.i.d.* standard normal random variables. Then define

$$s^*(d) = \sum_{j=1}^k \sum_{h \in D(t_j)} Z_h \left(\mathbf{x}_h - m_j n_j^{-1} \sum_{h \in R(t_j)} \mathbf{x}_h \right), \tag{11}$$

and

$$w^* = \sup_d s^{*'}(d) \tilde{v}^{-1}(d) s^*(d). \tag{12}$$

In Appendix C, we can show that the unconditional distribution of $\frac{1}{\sqrt{n}} \sup_d w(d)$ can be approximated by the conditional distribution of $\frac{1}{\sqrt{n}} w^*$. To this end, we generate the standard normal random sample (Z_1, \dots, Z_n) a large number of times. For each sample, we calculate w^* . The $100(1-\alpha)$ th percentile of the simulated w^* 's is the threshold value for the genome-wide significance level of α .

4. A SIMULATION STUDY

To investigate the proposed method in practical situations, we performed a small simulation study. Since the proposed score test is locally most powerful, we did not have to evaluate its power. In this section, we only examined the performance of the proposed interval-mapping method for locating the QTL for two different settings. The first setting was the same as the one in Diao [4] in order to compare their method, where the failure times were generated from

Table I. Sample means and standard errors for QTL location.

No. of markers	Weibull distribution		Log-normal distribution	
	Mean	SE	Mean	SE
6	35.6	10.7	33.4	10.0
11	33.4	9.8	33.9	8.4
51	33.0	8.9	32.3	6.3
101	32.1	7.5	33.6	7.9

the Weibull distribution with baseline hazard function $\lambda_0(t) = \gamma_1\gamma_2t^{\gamma_2-1}$ with $\gamma_1 = 0.01$ and $\gamma_2 = 2$. In the second setting, the failure time were generated from the log-normal distribution, that is, $\log T \sim N(0, 1)$ under the null hypothesis. In both settings, the censoring times were generated from the uniform $(0, \tau)$ distribution, where τ was chosen to yield $\sim 30\%$ censored observations. Assuming no crossover interference, we generated the marker data from the Markov chain. The interval-mapping step size was set at 1 cM.

We considered a chromosome with a total length of 100 cM. Genetic maps with different numbers of equally spaced markers were simulated. In both settings, one QTL located at 33 cM was simulated with $\beta_1 = 0.5$ and $\beta_2 = 0.25$. We generated 200 replicates of 300 observations from an F_2 population. The results are summarized in Table I where the unit of means and standard errors is cM.

In Table I, the results from setting 1 are very similar to those in Table 2 of Diao *et al.* [4]. In both settings, there is little bias for the estimation of the QTL location. In addition, we also found that the dense-maker map makes a small contribution to the accuracy of the confidence intervals of the QTL location.

ACKNOWLEDGEMENTS

I would like to thank Professor Daniel Rabinowitz for his help. And I am grateful to the editor and two referees for their helpful comments.

REFERENCES

- [1] Breslow N.E., Covariate analysis of censored survival data, *Biometrics* 30 (1974) 89–99.
- [2] Broman K.W., Mapping quantitative trait loci in the case of a spike in the phenotype distribution, *Genetics* 163 (2003) 1169–1175.
- [3] Churchill G.A., Doerge R.W., Empirical threshold values for quantitative trait mapping, *Genetics* 138 (1994) 963–971.

- [4] Diao G.Q., Lin D.Y., Zou F., Mapping quantitative trait loci with censored observations, *Genetics* 168 (2004) 1689–1698.
- [5] Dupuis J., Siegmund D., Statistical methods for mapping quantitative trait loci from a dense set of markers, *Genetics* 151 (1999) 373–386.
- [6] Fine J.P., Zou F., Yandell B.S., Nonparametric estimation of quantitative trait loci, *Biostatistics* 5 (2004) 501–513.
- [7] Haley C.S., Knott S.A., A simple regression method for mapping quantitative trait loci in line crosses using flanking markers, *Heredity* 69 (1992) 315–324.
- [8] Jansen R.C., Interval mapping of multiple quantitative trait loci, *Genetics* 135 (1993) 521–529.
- [9] Kalbfleisch J.D., Prentice R.L., *The Statistical Analysis of Failure Time Data*, 2nd edn., Wiley, Hoboken, NJ, 1998.
- [10] Kruglyak L., Lander E.S., A nonparametric approach for mapping quantitative trait loci, *Genetics* 139 (1995) 1421–1428.
- [11] Lander E.S., Botstein D., Mapping Mendelian factors underlying quantitative traits using RFLP linkage maps, *Genetics* 121 (1989) 185–199.
- [12] Lin D.Y., Wei L.J., Ying Z., Checking the Cox model with cumulative sums of martingale-based residuals, *Biometrika* 80 (1993) 557–572.
- [13] Lipsitz S.R., Ibrahim J.G., Estimating equations with incomplete categorical covariates in the Cox model, *Biometrics* 54 (1998) 1002–1023.
- [14] Lynch M., Walsh B., *Genetics and Analysis of Quantitative Traits*, Sinauer, Sunderland, MA, 1998.
- [15] Moreno C.R., Elsen J.M., Leroy P., Ducrocq V., Interval mapping methods for detecting QTL affecting survival and time-to-event phenotypes, *Genet. Res.* 85 (2005) 139–149.
- [16] Prentice R.L., Covariate measurement errors and parameter estimation in a failure time regression model, *Biometrika* 69 (1982) 331–342.
- [17] Symons R.C., Daly M.J., Fridlyand J., Speed T.P., Cook W.D., Gerondakis S., Harris A.W., Foote S.J., Multiple genetic loci modify susceptibility to plasmacytoma-related morbidity in $E\mu$ - v - abl transgenic mice, *Proc. Natl. Acad. Sci.* 99 (2002) 11299–11304.
- [18] Zeng Z.B., Precision mapping of quantitative traits loci, *Genetics* 136 (1994) 1457–1468.
- [19] Zou F., Fine J.P., Yandell B.S., On empirical likelihood for a semiparametric mixture model, *Biometrika* 89 (2002) 61–75.
- [20] Zou F., Fine J.P., Hu J., Lin D.Y., An efficient resampling method for assessing genome-wide statistical significance in mapping quantitative trait loci, *Genetics* 168 (2004) 2307–2316.

APPENDIX A: ASYMPTOTIC PROPERTIES

In this Appendix, we give some asymptotic results of model (2). The following arguments are similar to the Section 5.7 in Kalbfleisch and Prentice [9].

Define counting processes $\{N_i(t) = \Delta_i I(T_i \leq t), 0 \leq t\}$, at-risk processes $\{Y_i(t) = I(T_i \geq t, C_i \geq t), 0 \leq t\}, i = 1 \cdots, n$ and their generating filtration

$$\mathcal{F}_t = \sigma\{N_i(u), Y_i(u+), \mathbf{M}_i, i = 1, \cdots, n, 0 \leq u \leq t\}.$$

Then corresponding to this filtration, for each $i = 1, \cdots, n$, the process

$$M_i(t) = N_i(t) - \int_0^t Y_i(u)\lambda(u|\mathbf{M}_i)du, t \geq 0$$

is a martingale. Based on these martingales, the score of (5) can be written as

$$S(\boldsymbol{\beta}) = \sum_{i=1}^n \int \left(\frac{\mathbf{b}_i(t, \boldsymbol{\beta})}{a_i(t, \boldsymbol{\beta})} - \frac{\sum_{j=1}^n \mathbf{b}_j(t, \boldsymbol{\beta})Y_j(t)}{\sum_{j=1}^n a_j(t, \boldsymbol{\beta})Y_j(t)} \right) dM_i(t). \quad (\text{A.1})$$

This is equivalent to (6). Its predictable variation process is

$$\langle S(\boldsymbol{\beta}) \rangle = \sum_{i=1}^n \int \left(\frac{\mathbf{b}_i(t, \boldsymbol{\beta})}{a_i(t, \boldsymbol{\beta})} - \frac{\sum_{j=1}^n \mathbf{b}_j(t, \boldsymbol{\beta})Y_j(t)}{\sum_{j=1}^n a_j(t, \boldsymbol{\beta})Y_j(t)} \right)^{\otimes 2} a_i(t, \boldsymbol{\beta})Y_i(t)\lambda_0(t)dt, \quad (\text{A.2})$$

and an estimate of it at $\boldsymbol{\beta} = \mathbf{0}$ is

$$\sum_{i=1}^n \int \left(\frac{\mathbf{b}_i(t, \mathbf{0})}{a_i(t, \mathbf{0})} - \frac{\sum_{j=1}^n \mathbf{b}_j(t, \mathbf{0})Y_j(t)}{\sum_{j=1}^n a_j(t, \mathbf{0})Y_j(t)} \right)^{\otimes 2} \frac{a_i(t, \mathbf{0})Y_i(t)}{\sum_{j=1}^n a_j(t, \mathbf{0})Y_j(t)} dN_i(t), \quad (\text{A.3})$$

where $N_i(t) = \sum_{i=1}^n N_i(t)$. Since it is easy to verify $a_i(t, \mathbf{0}) = 1$ and $\mathbf{b}_i(t, \mathbf{0}) = \mathbf{x}_i$, for $i = 1, \cdots, n$, the above estimate is equivalent to (7).

APPENDIX B: ANALYTICAL APPROXIMATIONS OF THRESHOLDS

Let $\mathbf{g}_i(d) = (G_i(d), 1 - |G_i(d)|)'$, $i = 1, \cdots, n$, and $\mathbf{g} \stackrel{d}{=} \mathbf{g}_1$, where $G_i(d)$ is the genotype at position d . Let d_1 and d_2 denote two points on the chromosome, and p be the recombination fraction corresponding to the genetic distance $|d_1 - d_2|$. It is easy to see that the correlation of $\mathbf{g}(d_1)$ and $\mathbf{g}(d_2)$ is

$$\text{Corr}(\mathbf{g}(d_1), \mathbf{g}(d_2)) = \begin{pmatrix} 1 - 2p & 0 \\ 0 & 1 - 4p \end{pmatrix} = \begin{pmatrix} e^{-2|d_1 - d_2|} & 0 \\ 0 & e^{-4|d_1 - d_2|} \end{pmatrix},$$

assuming Haldane's map function.

The score statistic at $\beta = 0$ can be written as $s(d) = \sum_{i=1}^n \delta_i \{g_i(d) - \bar{g}(d, T_i)\}$, where $\bar{g}(d, t) = \sum_{i=1}^n g_i(d) Y_i(t) / \sum_{i=1}^n Y_i(t)$. Obviously, $\lim_{n \rightarrow \infty} \bar{g}(d, T_1) = E(g(d)) = (0, 1/2)'$, in probability. Define $u(d) = \tilde{v}(d)^{-1/2} s(d)$, we have

$$\text{Corr}(u(d_1), u(d_2)) = \text{Corr}(g(d_1), g(d_2)) + o_p(1).$$

Therefore, $u_1(d)$ and $u_2(d)$, two components of $u(d)$, are approximately independent Ornstein-Uhlenbeck processes with means zero and correlation functions $e^{-2|d_1-d_2|}$ and $e^{-4|d_1-d_2|}$, respectively. Then we can obtain analytical approximations of thresholds for both dense-map case and sparse-map case. See, for example, Dupuis and Siegmund [5].

APPENDIX C: RESAMPLING METHOD

It is easy to see that $\{s(d), d \geq 0\}$, where

$$s(d) = \sum_{i=1}^n \int \left(x_i - \frac{\sum_{j=1}^n x_j Y_j(t)}{\sum_{j=1}^n Y_j(t)} \right) dM_i(t), \quad (\text{C.1})$$

converges to a zero-mean Gaussian process. We can approximate its limiting distribution through the Monte Carlo method. A robust and efficient method is to replace $M_i(t)$ by a similar process, say $\tilde{M}_i(t)$, which has a known distribution and leave other terms unchanged. Note that the variance function of $M_i(t)$ is $E[N_i(t)]$. Thus a natural candidate for $\tilde{M}_i(t)$ is $N_i(t)Z_i$, where $\{Z_i, i = 1 \cdots, n\}$ denotes a random sample of standard normal variables. After doing that, we obtain

$$s^*(d) = \sum_{i=1}^n \int Z_i \left(x_i - \frac{\sum_{j=1}^n x_j Y_j(t)}{\sum_{j=1}^n Y_j(t)} \right) dN_i(t). \quad (\text{C.2})$$

By the arguments similar to the Appendix 1 in Lin *et al.* [12], we can show that under the null hypothesis that there is no QTL, the conditional distribution of $n^{-\frac{1}{2}} s^*(d)$ given the observed data is the same in the limit as the unconditional distribution of $n^{-\frac{1}{2}} s(d)$.