



HAL
open science

Likelihood Ratio Test process for Quantitative Trait Locus detection

Charles-Elie Rabier, Jean-Marc Azaïs, Céline Delmas

► **To cite this version:**

Charles-Elie Rabier, Jean-Marc Azaïs, Céline Delmas. Likelihood Ratio Test process for Quantitative Trait Locus detection. 2010. hal-00483171v1

HAL Id: hal-00483171

<https://hal.science/hal-00483171v1>

Preprint submitted on 12 May 2010 (v1), last revised 17 Dec 2012 (v4)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Likelihood Ratio Test process for Quantitative Trait Locus detection

Charles-Elie Rabier

*Université de Toulouse, Institut de Mathématiques de Toulouse, U.P.S., Toulouse, France.
INRA UR631, Station d'Amélioration Génétique des Animaux, Auzeville, France.*

Jean-Marc Azaïs

Université de Toulouse, Institut de Mathématiques de Toulouse, U.P.S., Toulouse, France.

Céline Delmas

INRA UR631, Station d'Amélioration Génétique des Animaux, Auzeville, France.

Summary. We consider the likelihood ratio test (LRT) process related to the test of the absence of QTL on the interval $[0, T]$ representing a chromosome (a QTL denotes a quantitative trait locus, i.e. a gene with quantitative effect on a trait). We give the asymptotic distribution of this LRT process under the null hypothesis that there is no QTL on $[0, T]$ and under the alternative that there exists a QTL at t^* on $[0, T]$. We show that the LRT is asymptotically the square of a non linear interpolated process. We propose a simple and original method to calculate the maximum and the argmax of the LRT process using only statistics on markers and their ratio. We finally propose a new method to calculate thresholds for QTL detection.

Keywords: Gaussian process, Likelihood Ratio Test, Mixture models, Nuisance parameters present only under the alternative, QTL detection, χ^2 process.

1. Introduction

We study a backcross population: $A \times (A \times B)$, where A and B are purely homozygous lines and we address the problem of detecting a Quantitative Trait Locus, so-called QTL (a gene influencing a quantitative trait which is able to be measured) on a given chromosome. The trait is observed on n individuals (progenies) and we denote by Y_j , $j = 1, \dots, n$, the observations, which we will assume to be independent and identically distributed (iid). The mechanism of genetics, or more precisely of meiosis, implies that among the two chromosomes of each individual, one is purely inherited from A while the other (the "recombined" one), consists of parts originated from A and parts originated from B , due to crossing-overs. Using the Haldane (1919) distance and modelling, each chromosome will be represented by a segment $[0, T]$. The distance on $[0, T]$ is called the genetic distance (which is measured in Morgans). The key point is that, if the true position of the QTL is $t = t^*$, the response Y obeys to a mixture model with known weights :

$$p(t)f_{(\mu+q,\sigma)}(\cdot) + \{1 - p(t)\} f_{(\mu-q,\sigma)}(\cdot) \quad (1)$$

where $f_{(\mu,\sigma)}(\cdot)$ denotes a Gaussian density with mean μ and variance σ^2 . (μ, q, σ) are the unknown parameters. At every location $t \in [0, T]$, we perform a likelihood ratio test (LRT) of the hypothesis " $q = 0$ " in formula (1) based on n observations Y_1, \dots, Y_n . We call $\Lambda_n(t)$ the obtained quantity. The dependence on t of the weights is precisely described

in Section 3. We denote $p_j(t)$ the value of the weight $p(t)$ for the j th observation. The process $\{\Lambda_n(t), t \in [0, T]\}$ will be called "likelihood ratio test process" and taking as test statistic the maximum of this process comes down to perform a LRT in a model when the localisation of the QTL is an extra parameter.

In the special case where the weights are 0 or 1 depending on the individual, Lander and Botstein (1989) stated that the asymptotic distribution of the LRT process along $[0, T]$ is the square of an Ornstein-Uhlenbeck process. This result has been proved by Cierco (1998). Bounds for the distribution of the maximum of a regularization of an Ornstein-Uhlenbeck process were proposed by Azaïs and Cierco-Ayrolles (2002), Azaïs and Wschebor (2009). Some results about the asymptotic distribution of the LRT process under the null hypothesis are given in Rebaï et al. (1994) for a special modelling of the weights. Their results are inferred from the bounds given by Davies (1977), Davies (1987) for the maximum of sufficiently regular Gaussian and chi-square processes.

In this paper we consider the modelling of the weight used by geneticists to detect QTL, so called Interval Mapping. We give the asymptotic distribution of the LRT process along the interval $[0, T]$ under the null hypothesis that there is no QTL on $[0, T]$ ($q = 0$) and under the alternative that there is one QTL at t^* on $[0, T]$ which means that the quantitative trait for each individual is distributed as the mixture in formula (1) with $t = t^*$.

The main result of the paper (theorem 1 and theorem 3) is that the LRT process is asymptotically the square of a "non linear interpolated process". It describes the fact that, when we analyze data, the likelihood profile (ie. the path observed of the LRT process) is smooth between markers. Besides, we have a close formula (lemma 1 and lemma 2) to compute the maximum of the LRT process. This formula allows us to give advice on how to analyze data : we should first perform tests on markers and then calculate only one other statistic in each marker interval if the ratio between the score statistics on the flanking markers fulfill a given condition. Finally, we propose a new method suitable whatever the genetic map, using Monte-Carlo Quasi Monte-Carlo (Genz (1992)), to calculate thresholds for QTL detection. This method will be compared with Rebaï et al. (1994)'s method based on Davies (1977), and with Feingold and al. (1993)'s method based on Siegmund (1985).

Note that in this article, we also prove that the LRT process obtained by Rebaï et al. (1994), Rebaï et al. (1995) is asymptotically the square of a "linear interpolated process" and we generalize their results to the alternative hypothesis. Besides, we show that the law of the maximum of the square of the "non linear interpolated process" is the same as the law of the maximum of the square of the "linear interpolated process". We refer to the book of Van der Vaart (1998) for element of asymptotic statistics used in proofs.

2. Model

The chromosome is the segment $[0, T]$. K genetic markers are located on the chromosome, one at each extremity. $t_1 = 0 < t_2 < \dots < t_K = T$ are the locations of the markers. The "genome information" at t will be denoted $X(t)$. The Haldane (1919) model can be written mathematically : let $N(t)$ be a standard Poisson process, the law of $X(t)$ is $\frac{1}{2}(\delta_1 + \delta_{-1})$ and $X(t) = (-1)^{N(t)}X(t_1)$. The Haldane (1919)'s function $r : [0, T]^2 \mapsto [0, \frac{1}{2}]$ is such as

$$r(t, t') = \mathbb{P}(X(t)X(t') = -1) = \mathbb{P}(|N(t) - N(t')| \text{ odd}) = \frac{1}{2} (1 - e^{-2|t-t'|})$$

$\bar{r}(t, t')$ will be the function equal to $1 - r(t, t')$.

We are interested in a quantitative trait Y which depends on the value of $X(t)$ at $t^* \in [t_1, t_K]$ which is the location of the QTL. The quantitative trait verifies :

$$Y_j = \mu + X(t^*) q + \sigma \varepsilon$$

where ε is a Gaussian white noise and q the effect of the QTL.

Besides, the "genome information" is available only at locations of genetic markers, that is to say at t_1, t_2, \dots, t_K . We denote by $X_j(t)$ the value of the variable $X(t)$ for the j th observation. So, in fact, our observation on each individual is $(Y_j, X_j(t_1), \dots, X_j(t_K))$. These observations are supposed to be iid. The goal of this study is to test if q is equal to zero. The challenge is that t^* is unknown.

3. Only 2 genetic markers

To begin, we suppose that there are only two markers ($K = 2$) located at 0 and $T : 0 = t_1 < t_2 = T$. As explained previously, we are looking for a QTL lying at a position $t^* \in [t_1, t_2]$. Let $t \in [t_1, t_2]$. It is clear that the weight $p(t)$ satisfies $p(t) = \mathbb{P}\{X(t) = 1 | X(t_1), X(t_2)\}$. Consider for example the case $X(t_1) = X(t_2) = 1$, then by the Bayes rule :

$$\mathbb{P}\{X(t) = 1 | X(t_1) = 1, X(t_2) = 1\} = \frac{(1/2) \mathbb{P}\{N(t) - N(t_1) \text{ even}\} \mathbb{P}\{N(t_2) - N(t) \text{ even}\}}{(1/2) \mathbb{P}\{N(t_2) - N(t_1) \text{ even}\}}$$

So that, in general $\forall t \in]t_1, t_2[$:

$$p(t) = Q_t^{1,1} 1_{X(t_1)=1} 1_{X(t_2)=1} + Q_t^{1,-1} 1_{X(t_1)=1} 1_{X(t_2)=-1} + Q_t^{-1,1} 1_{X(t_1)=-1} 1_{X(t_2)=1} + Q_t^{-1,-1} 1_{X(t_1)=-1} 1_{X(t_2)=-1} \quad (2)$$

where :

$$Q_t^{1,1} = \frac{\bar{r}(t_1, t) \bar{r}(t, t_2)}{\bar{r}(t_1, t_2)}, \quad Q_t^{1,-1} = \frac{\bar{r}(t_1, t) r(t, t_2)}{r(t_1, t_2)}$$

$$Q_t^{-1,1} = \frac{r(t_1, t) \bar{r}(t, t_2)}{r(t_1, t_2)}, \quad Q_t^{-1,-1} = \frac{r(t_1, t) r(t, t_2)}{\bar{r}(t_1, t_2)}$$

We can remark that we have :

$$Q_t^{-1,-1} = 1 - Q_t^{1,1} \quad \text{and} \quad Q_t^{-1,1} = 1 - Q_t^{1,-1}$$

Besides, $p(t_1) = 1_{X(t_1)=1}$ and $p(t_2) = 1_{X(t_2)=1}$. So, the weights $p(t)$ are continuous at t_1 and t_2 .

Let $\theta = (q, \mu, \sigma)$ be the parameter of the model at t fixed and $\theta_0 = (0, \mu, \sigma)$ the true value of the parameter under H_0 . The likelihood of the triplet $(Y, X(t_1), X(t_2))$ with respect to the measure $\lambda \otimes N \otimes N$, λ being the Lebesgue measure, N the county measure on \mathbb{N} , is $\forall t \in [t_1, t_2]$:

$$L(\theta, t) = [p(t) f_{(\mu+q, \sigma)}(y) + \{1 - p(t)\} f_{(\mu-q, \sigma)}(y)] g(t) \quad (3)$$

4 Céline Delmas

where

$$g(t) = \frac{1}{2} \{ \bar{r}(t_1, t_2) 1_{X(t_1)=1} 1_{X(t_2)=1} + r(t_1, t_2) 1_{X(t_1)=1} 1_{X(t_2)=-1} \} \\ + \frac{1}{2} \{ r(t_1, t_2) 1_{X(t_1)=-1} 1_{X(t_2)=1} + \bar{r}(t_1, t_2) 1_{X(t_1)=-1} 1_{X(t_2)=-1} \}$$

The likelihood $L_n(\theta, t)$ for n observations is obtained by the product of n terms as above. $\hat{\theta} = (\hat{q}, \hat{\mu}, \hat{\sigma})$ will be the maximum likelihood estimator (MLE) of θ .

Under H_0 , there is no QTL lying on the interval $[t_1, t_2]$. Besides, under H_1 , it is supposed that there is only one location where the QTL lies. The location of the QTL, t^* ($t^* \in [t_1, t_2]$), will be added in the definition of H_1 . So, the alternative hypothesis can be written :

$$H_{at^*} : \text{"the QTL is located at the position } t^* \text{ with effect } q = a/\sqrt{n} \text{ where } a \in \mathbb{R}^* \text{"}$$

The QTL effect q is such as $q = a/\sqrt{n}$ in order to deal with Le Cam (1986)'s theory.

3.1. A "non linear interpolated process"

Theorem 1 *With the previous defined notations, and defining respectively $\Lambda_n(\cdot)$ and $S_n(\cdot)$, the LRT process and the score process for n observations,*

$$S_n(\cdot) \Rightarrow Z(\cdot) \quad , \quad \Lambda_n(\cdot) \xrightarrow{F.d.} \{Z(\cdot)\}^2$$

as n tends to infinity, under H_0 and H_{at^*} where :

- \Rightarrow is the weak convergence and $\xrightarrow{F.d.}$ is the convergence of finite-dimensional distributions
- $Z(\cdot)$ is the Gaussian process with covariance function $\forall(t, t') \in [t_1, t_2]^2$:

$$\Gamma(t, t') = \frac{4\mathbb{E}\{p(t)p(t')\} - 1}{\sqrt{\mathbb{E}\left[\{2p(t) - 1\}^2\right]} \sqrt{\mathbb{E}\left[\{2p(t') - 1\}^2\right]}}$$

and expectation $\forall(t, t^*) \in [t_1, t_2]^2$:

- under H_0 , $m(t) = 0$
- under H_{at^*}

$$m_{t^*}(t) = \frac{a \mathbb{E}[X(t^*) \{2p(t) - 1\}]}{\sigma \sqrt{\mathbb{E}\left[\{2p(t) - 1\}^2\right]}}$$

Another way of characterizing $Z(\cdot)$ is that $Z(\cdot)$ is the non linear interpolated process such as $\forall t \in [t_1, t_2]$:

$$Z(t) = \{ \alpha(t) Z(t_1) + \beta(t) Z(t_2) \} / \sqrt{\mathbb{E}\left[\{2p(t) - 1\}^2\right]}$$

where $\forall t \in]t_1, t_2[$, $\alpha(t) = Q_t^{1,1} + Q_t^{1,-1} - 1$, $\beta(t) = Q_t^{1,1} - Q_t^{1,-1}$ and $\alpha(t_1) = 1$, $\beta(t_1) = 0$, $\alpha(t_2) = 0$, $\beta(t_2) = 1$, $\text{Cov}\{Z(t_1), Z(t_2)\} = e^{-2(t_2-t_1)}$.

In the same way, $\forall (t, t^*) \in [t_1, t_2]^2$:

$$m_{t^*}(t) = \{ \alpha(t) m_{t^*}(t_1) + \beta(t) m_{t^*}(t_2) \} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

The quantity $\mathbb{E} \left[\{2p(t) - 1\}^2 \right]$ is given in formula (9) of the proof of the theorem in Section 5.1. $\mathbb{E} \{p(t)p(t')\}$ is given in appendix 7.1. $\mathbb{E} [X(t^*) \{2p(t) - 1\}]$ is given in formula (14) of the proof in Section 5.1.

We limit our attention to finite dimensional convergence since for the applications, the interval studied is always discretized, Wu et al. (2007).

Figures 1 represent the covariance function $\Gamma(t, t')$ and also the mean function $m_{t^*}(t)$. T is equal to 0.2M. We can remark that the covariance function is regular.

Contrary to Azaïs et al. (2006) and Azaïs et al. (2009), the shift at position t is not $\Gamma(t, t^*)$. The model considered here is more complicated due to the fact that an observation includes the quantitative trait Y and the "genome information", $X(t_1)$ and $X(t_2)$.

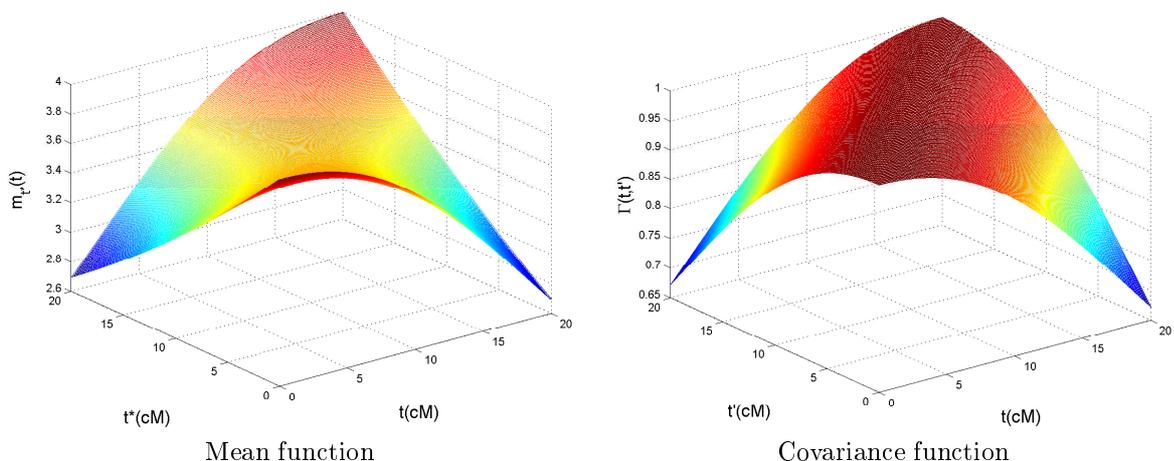


Fig. 1. Mean function and Covariance function ($\alpha = 4$, $\sigma = 1$, $T = 0.2M$)

3.2. Remarks

As it is well known, for regular model, LRT is equivalent to score test in the sense that $\forall t \in [t_1, t_2]$:

$$\Lambda_n(t) = \{S_n(t)\}^2 + o_{P_{\theta_0}}(1)$$

We remind that, as in the proof of the theorem in Section 5.1, the notation $o_{P_{\theta_0}}(1)$ is short for a sequence of random vectors that converges to zero in probability under H_0 (i.e. no

QTL on the whole interval studied).

A little algebra shows (see Section 5.1) :

$$S_n(t) = \sum_{j=1}^n \frac{(y_j - \mu) (2p_j(t) - 1)}{\sqrt{n} \sigma \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}} \quad (4)$$

The score test can be obtained, replacing μ by $\hat{\mu} = \sum y_j/n$, according to Prohorov, and replacing σ by $\hat{\sigma} = \sqrt{\sum (y_j - \hat{\mu})^2/n}$, according to Slutsky's lemma. Nevertheless, in this article, in order to make the reading easier, the score test statistic is defined as in formula (4). The score process considered in theorem 1 is based on this formula. However, we have the same result as in theorem 1 for the other score process because the tightness of this process is obvious according to the proof of theorem 1.

After some calculations, we can remark that :

$$S_n(t) = \{ \alpha(t) S_n(t_1) + \beta(t) S_n(t_2) \} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]} \quad (5)$$

with $\text{Cov}_{H_0} \{S_n(t_1), S_n(t_2)\} = e^{-2(t_2-t_1)}$. $\alpha(t)$ and $\beta(t)$ are given quantities in theorem 1. It comes :

$$\Lambda_n(t) = \{ \alpha(t) S_n(t_1) + \beta(t) S_n(t_2) \}^2 / \mathbb{E} \left[\{2p(t) - 1\}^2 \right] + o_{P_{\theta_0}}(1)$$

Besides, by contiguity (cf. proof of theorem 1 in Section 5.1), the quantity $o_{P_{\theta_0}}(1)$ converges also to zero under H_{at^*} . That is to say, the LRT statistic at a position t between the two genetic markers is asymptotically equal to the square of a non linear interpolation between the score test statistics on the markers.

3.3. A "linear interpolated process"

To construct an approximation of $S_n(\cdot)$ (and $\Lambda_n(\cdot)$), we introduce a new process $V_n(\cdot)$ which is obtained from $S_n(\cdot)$ by :

- linear (or polygonal) interpolation
- renormalization

More precisely :

$$V_n(t) = \left\{ \frac{t_2 - t}{t_2 - t_1} S_n(t_1) + \frac{t - t_1}{t_2 - t_1} S_n(t_2) \right\} / \sqrt{\tau(t)} \quad (6)$$

where

$$\begin{aligned} \tau(t) &= \mathbb{V}_{H_0} \left\{ \frac{t_2 - t}{t_2 - t_1} S_n(t_1) + \frac{t - t_1}{t_2 - t_1} S_n(t_2) \right\} \\ &= \left(\frac{t_2 - t}{t_2 - t_1} \right)^2 + 2 \frac{(t - t_1)(t_2 - t)}{(t_2 - t_1)^2} e^{-2(t_2 - t_1)} + \left(\frac{t - t_1}{t_2 - t_1} \right)^2 \end{aligned}$$

It can be seen easily that $\tau(t) \neq 0$, $\forall t \in [t_1, t_2]$. $V_n(\cdot)$ remains asymptotically a Gaussian process, centered under H_0 , with unit variance and $\text{Cov}_{H_0} \{S_n(t_1), S_n(t_2)\} = e^{-2(t_2-t_1)}$.

Some comments about the process $V_n(\cdot)$:

- (a) According to formula (10) in Section 5.1 and after some calculations, we can establish that asymptotically, the process $V_n^2(\cdot)$ corresponds to likelihood ratio tests for a mixture model whose weights verify :

$$p(t) = 1_{X(t_1)=1}1_{X(t_2)=1} + \frac{t_2 - t}{t_2 - t_1} 1_{X(t_1)=1}1_{X(t_2)=-1} + \frac{t - t_1}{t_2 - t_1} 1_{X(t_1)=-1}1_{X(t_2)=1} \quad (7)$$

We can remark that these weights are an approximation at the first order of the weights considered previously in formula (2). So, $V_n(\cdot)$ will be a good approximation if and only if the genetic markers are close to each other.

- (b) $V_n^2(\cdot)$ is a generalization of the process studied, under H_0 , by Rebaï et al. (1995) : the number of individuals in each class is not equal to the expectations (respectively $n\bar{r}(t_1, t_2)/2$, $nr(t_1, t_2)/2$, $nr(t_1, t_2)/2$, $n\bar{r}(t_1, t_2)/2$) but is still random (respectively $\sum_{j=1}^n 1_{X_j(t_1)=1}1_{X_j(t_2)=1}$, $\sum_{j=1}^n 1_{X_j(t_1)=1}1_{X_j(t_2)=-1}$, $\sum_{j=1}^n 1_{X_j(t_1)=-1}1_{X_j(t_2)=1}$ and $\sum_{j=1}^n 1_{X_j(t_1)=-1}1_{X_j(t_2)=-1}$).
- (c) By contiguity (cf. proof of theorem 1 in Section 5.1), under H_{at^*} , $V_n(\cdot)$ is asymptotically the same process as under H_0 on which the mean function $\tilde{m}_{t^*}(t)$ has been added. $\tilde{m}_{t^*}(t)$ is such as :

$$\tilde{m}_{t^*}(t) = \left\{ \frac{t_2 - t}{t_2 - t_1} m_{t^*}(t_1) + \frac{t - t_1}{t_2 - t_1} m_{t^*}(t_2) \right\} / \sqrt{\tau(t)}$$

- (d) $V_n(\cdot)$ is defined here with $\text{Cov}_{H_0} \{S_n(t_1), S_n(t_2)\} = e^{-2(t_2-t_1)}$. In order to consider other covariances between $S_n(t_1)$ and $S_n(t_2)$, $\tau(\cdot)$ has to be adapted. It can easily be seen that the new process $V_n^2(\cdot)$ is still a generalization of the process studied by Rebaï et al. (1995) provided that $\mathbb{E} \left[\{2p(t) - 1\}^2 \right] \neq 0$ ($p(t)$ verifies formula (7)).

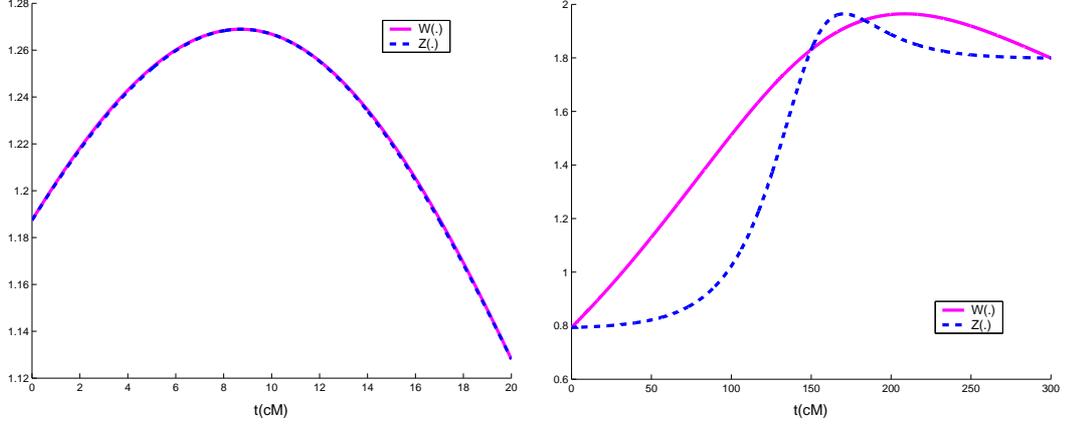
We will name $W(\cdot)$ the limiting process of $V_n(\cdot)$: $W(\cdot)$ is the linear interpolated process. Figure 2 represents two paths of the processes $W(\cdot)$ and $Z(\cdot)$ under the null hypothesis. We remind that $Z(\cdot)$ is the non linear interpolated process. We can observe that the paths of the two processes overlap when the distance between the two genetic markers is 20cM, that is to say 0.2M. As mentioned before, $W(\cdot)$ is a good approximation of $Z(\cdot)$ when the genetic markers are close to each other. As expected, we can remark that when the distance between the markers is 3M, the paths of the two processes don't overlap anymore. Note that same conclusions hold under contiguous alternatives (cf. Rabier (2010)).

3.4. Impact of the interpolations on data analysis

In this Section, we state why the results about the non linear interpolation and about the linear interpolation are important for data analysis. To begin, we present a theorem and a lemma.

Theorem 2 Let C_1 and C_2 be two continuous random variables, and let $\tilde{\rho}$ such as $0 < \tilde{\rho} < 1$. Let consider $\gamma_1(t)$ and $\gamma_2(t)$, two continuous functions on $[t_1, t_2]$, and let define the process $D(\cdot)$ on $[t_1, t_2]$ such as :

$$D(t) = \frac{\gamma_1(t) C_1 + \gamma_2(t) C_2}{\sqrt{\{\gamma_1(t)\}^2 + \{\gamma_2(t)\}^2 + 2 \tilde{\rho} \gamma_1(t) \gamma_2(t)}}$$



One path for the two processes ($T = 0.2M$) One path for the two processes ($T = 3M$)

Fig. 2. Comparison between the paths, under H_0 , of the linear interpolated process $W_{(\cdot)}$ and those of the non linear interpolated process $Z_{(\cdot)}$.

then if the function $\frac{\gamma_2(t)}{\gamma_1(t) + \gamma_2(t)}$ is bounded by 0 and 1 on $[t_1, t_2]$, and if these bounds are reached, then

$$\sup_{t \in [t_1, t_2]} \{D(t)\}^2 = \max \left[\{C_1\}^2, \frac{\{C_1\}^2 + \{C_2\}^2 - 2 \tilde{\rho} C_1 C_2}{(1 + \tilde{\rho})(1 - \tilde{\rho})} \mathbf{1}_{\frac{C_2}{C_1} \in] \tilde{\rho}, \frac{1}{\tilde{\rho}} [}, \{C_2\}^2 \right]$$

Lemma 1 *With the previous defined notations and reminding that $W(t_1) = Z(t_1)$ and $W(t_2) = Z(t_2)$,*

let $\xi = \frac{(t_2 - t_1) \{e^{-2(t_2 - t_1)} W(t_1) - W(t_2)\}}{\{e^{-2(t_2 - t_1)} - 1\} \{W(t_1) + W(t_2)\}} + t_1$, then under H_0 and H_{at^*}

$$\{W(\xi)\}^2 = \frac{\{W(t_1)\}^2 + \{W(t_2)\}^2 - 2 e^{-2(t_2 - t_1)} W(t_1) W(t_2)}{\{1 + e^{-2(t_2 - t_1)}\} \{1 - e^{-2(t_2 - t_1)}\}} \quad \text{and}$$

$$\begin{aligned} \sup_{t \in [t_1, t_2]} \{Z(t)\}^2 &= \sup_{t \in [t_1, t_2]} \{W(t)\}^2 \\ &= \max \left[\{W(t_1)\}^2, \{W(\xi)\}^2 \mathbf{1}_{\frac{W(t_2)}{W(t_1)} \in] e^{-2(t_2 - t_1)}, e^{2(t_2 - t_1)} [}, \{W(t_2)\}^2 \right] \end{aligned}$$

The proof of theorem 2 and lemma 1 are respectively given in Sections 5.2 and 5.3. According to lemma 1, even when the genetic markers are not close to each other, the law of the supremum of the square of the two interpolated processes is the same (see Figure 2). However, when the supremum is obtained between markers, it is not obtained at the same

positions. These locations are ξ for $\{W(\cdot)\}^2$ and ξ' for $\{Z(\cdot)\}^2$, where ξ' is such as :

$$\frac{(t_2 - t_1) \beta(\xi')}{\alpha(\xi') + \beta(\xi')} + t_1 = \xi$$

On the other hand, lemma 1 can easily be adapted to the non asymptotic processes. Indeed, we can replace $Z(\cdot)$ by $S_n(\cdot)$, and $W(\cdot)$ by $V_n(\cdot)$ everywhere in this lemma, since the focus is on the same interpolations. As previously, when the supremum is obtained between markers, it is obtained at ξ for $V_n(\cdot)$ and ξ' for $S_n(\cdot)$. Furthermore, we can advise not to perform a large number of tests between the genetic markers anymore. First, we should perform score tests on markers. Then, only if the ratio observed between the score statistics on markers (ie. the ratio $S_n(t_2)/S_n(t_1)$) belongs to the interval $] e^{-2(t_2-t_1)}, e^{2(t_2-t_1)} [$, we have to calculate another quantity :

$$\zeta = \frac{\{S_n(t_1)\}^2 + \{S_n(t_2)\}^2 - 2 e^{-2(t_2-t_1)} S_n(t_1) S_n(t_2)}{\{1 + e^{-2(t_2-t_1)}\} \{1 - e^{-2(t_2-t_1)}\}}$$

To conclude, we should use as a test statistic :

$$\max \left[\{S_n(t_1)\}^2, \zeta \mathbb{1}_{\frac{S_n(t_2)}{S_n(t_1)} \in] e^{-2(t_2-t_1)}, e^{2(t_2-t_1)} [}, \{S_n(t_2)\}^2 \right]$$

4. Several markers : the ‘‘Interval Mapping’’ of Lander and Botstein (1989)

In that case suppose that there are K markers $0 = t_1 < t_2 < \dots < t_K = T$. We consider values t, t' or t^* of the parameters that are distinct of the markers positions, and the result will be prolonged by continuity at the markers positions. For $t \in [t_1, t_K] \setminus \mathbb{T}_k$ where $\mathbb{T}_k = \{t_1, \dots, t_K\}$, we define t^ℓ and t^r as :

$$t^\ell = \sup \{t_k \in \mathbb{T}_k : t_k < t\}, \quad t^r = \inf \{t_k \in \mathbb{T}_k : t < t_k\}$$

In other words, t belongs to the ‘‘Marker interval’’ (t^ℓ, t^r).

Theorem 3 *We have the same result as in theorem 1 except that the following expressions are more complicated :*

$$\mathbb{E} \left[\{2p(t) - 1\}^2 \right], \quad \mathbb{E} \{p(t)p(t')\}, \quad \mathbb{E} [X(t^*) \{2p(t) - 1\}], \quad \alpha(t), \quad \beta(t)$$

Besides, $Z(\cdot)$ is now the non linear interpolated process such as :

$$Z(t) = \{ \alpha(t) Z(t^\ell) + \beta(t) Z(t^r) \} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

with $\forall k \forall k', \text{Cov} \{Z(t_k), Z(t_{k'})\} = e^{-2|t_k - t_{k'}|}$.

In the same way, the mean function $m_{t^*}(t)$ is now such as :

$$m_{t^*}(t) = \{ \alpha(t) m_{t^*}(t^\ell) + \beta(t) m_{t^*}(t^r) \} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

All these expressions including a proof are given in appendix 7.2.

Note that $\forall k \ \forall k', \Gamma(t_k, t_{k'}) = e^{-2|t_k - t_{k'}|}$. It is relative to an Ornstein-Uhlenbeck process, as studied in Lander and Botstein (1989), and Cierco (1998).

The paths of three processes are presented in Figure 3 ($T = 1M$):

- the Ornstein-Uhlenbeck process.
- the process $Z(\cdot)$ with only 2 markers, located at $t_1 = 0$ and $t_2 = 1M$.
- the process $Z(\cdot)$ with markers located every 10cM.

The paths of the last two processes are smooth (due to the interpolation) whereas the paths of the Ornstein-Uhlenbeck process are very jerky. It's not surprising because the Ornstein-Uhlenbeck process can be viewed as a stationary version of the Brownian motion.

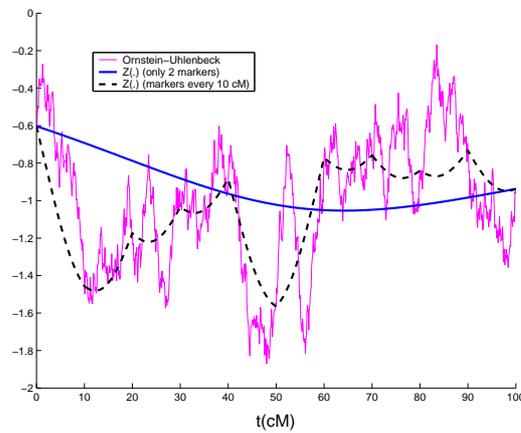


Fig. 3. Paths of three different Gaussian processes

Recently, the law of the LRT process under the null hypothesis has also been obtained by Chang et al. (2009). Technical differences are presented in appendix 7.4. The originality of our work is first, that we focus not only on the null hypothesis. Secondly, we show that the LRT process is asymptotically the square of a “non linear interpolated process”. It describes the fact that, when we analyze data, the likelihood profile (ie. the path observed of the LRT process) is smooth between markers.

4.1. Remarks

The linear interpolated process $W(\cdot)$ presented in Section 3.3 can easily be generalized to the case of several markers. This is a generalization of the process studied, under H_0 , by Rebaï et al. (1994). The details are given in appendix 7.3.

In the same way, lemma 1 can be generalized :

Lemma 2 *With the previous defined notations*

and reminding that $\forall k \ W(t_k) = Z(t_k)$,

$$\text{let } \xi(t^\ell, t^r) = \frac{(t^r - t^\ell) \left\{ e^{-2(t^r - t^\ell)} W(t^\ell) - W(t^r) \right\}}{\left\{ e^{-2(t^r - t^\ell)} - 1 \right\} \left\{ W(t^\ell) + W(t^r) \right\}} + t^\ell$$

, then under H_0 and H_{at^*}

$$[W \{ \xi(t^\ell, t^r) \}]^2 = \frac{\{W(t^\ell)\}^2 + \{W(t^r)\}^2 - 2 e^{-2(t^r - t^\ell)} W(t^\ell) W(t^r)}{\{1 + e^{-2(t^r - t^\ell)}\} \{1 - e^{-2(t^r - t^\ell)}\}}$$

and

$$\sup_{t \in [0, T]} \{Z(t)\}^2 = \sup_{t \in [0, T]} \{W(t)\}^2$$

$$\sup_{t \in [0, T]} \{W(t)\}^2 = \max \left[\{W(t_1)\}^2, \dots, \{W(t_K)\}^2, [W \{ \xi(t_1, t_2) \}]^2 \mathbf{1}_{\frac{W(t_2)}{W(t_1)} \in] e^{-2(t_2 - t_1)}, e^{2(t_2 - t_1)} [}, \dots, [W \{ \xi(t_{K-1}, t_K) \}]^2 \mathbf{1}_{\frac{W(t_K)}{W(t_{K-1})} \in] e^{-2(t_K - t_{K-1})}, e^{2(t_K - t_{K-1})} [} \right]$$

The proof of this lemma is largely inspired from the proof of lemma 1. Using lemma 2 and using the same arguments as in Section 3.4, we can advise to perform tests on markers and to calculate one other statistic in each marker interval when it is required.

4.2. Application to the calculation of thresholds

The theoretical results presented in this article allow us to propose a new method to obtain the $\alpha\%$ quantile of the supremum of the process $\{Z(\cdot)\}^2$ under H_0 . This method is a direct application of lemma 2. Besides, Monte-Carlo Quasi Monte-Carlo (MCQMC) methods of Genz (1992) which are very fast have been considered. As the numerical computation of a multivariate normal distribution is often a difficult problem, Genz described in his paper, a transformation that simplifies the problem and places into a form that allows efficient calculations using standard numerical multiple integration algorithms. He suggests to use in particular MCQMC algorithms. Indeed, a simple Monte-Carlo method (MC) using N points have errors that are typically $O(1/\sqrt{N})$ whereas Quasi Monte-Carlo methods (QMC) have errors $O(1/N)$. In order to be sure that the functions studied have nice properties for QMC, another Monte-Carlo step is required, this is MCQMC. We refer to Genz (1992) for more details. We use here function QSIMVNEF of Genz, which is a Matlab function with supporting functions, for the numerical computation of multivariate normal distribution expected values. This function has been adapted and a Newton method has been used in order to obtain the thresholds.

Our method is available in a Matlab package with graphical user interface : "imapping.zip". It can be downloaded at www.math.univ-toulouse.fr/~rabier.

In this Section, we propose to compare the performances of our method with other methods usually used in QTL detection.

In Rebaï et al. (1994), we can find an upper bound for the threshold. This bound is the

quantity c^2 such as :

$$1 - \alpha = 2 \Phi(-c) + \frac{2 e^{-c^2/2}}{\pi} \sum_{k=1}^{K-1} \arctan \left(\sqrt{\frac{1 - e^{-2(t_{k+1}-t_k)}}{1 + e^{-2(t_{k+1}-t_k)}}} \right)$$

where Φ is the cumulative distribution of the standardized normal distribution. This method is based on Davies (1977). However, it is sensitive to the number of genetic markers. Indeed, the derivative of the process $W(\cdot)$ has a jump at each markers location, and Davies (1977) upper bound is suitable when the derivative of the process has a finite number of jumps.

In Feingold and al. (1993), the authors propose a threshold based on the discrete process resulting from tests only on markers. Besides, they suppose constant the distance between genetic markers. The threshold c^2 is such as :

$$1 - \alpha = 1 - \Phi(c) + 2 T c \varphi(c) \nu(2c\sqrt{\Delta})$$

where φ is the density of a normal standardized, Δ is distance between two consecutive markers.

This method is inspired from Siegmund (1985) where the function ν is fully described.

In Figures 4 and 5, thresholds corresponding to different methods are computed. As expected, Rebaï's method is very sensitive to the number of genetic markers. We can observe that Feingold's method and our method give almost same results.

However, our method give different results than Feingold when the number of genetic markers is very small (cf. Figure 6). Indeed, Feingold's method requires the number of genetic makers to be not too small (cf. Feingold and al. (1993)). The advantage of our method is that this method is appropriate whatever the map.

| Method | <i>this paper</i> | <i>Rebaï</i> | <i>Feingold</i> |
|-----------|-------------------|--------------|-----------------|
| Threshold | 6.76 | 6.92 | 6.78 |

Fig. 4. Thresholds as a function of the method considered. The map consists of 6 genetic markers equally spaced every 20cM ($T = 1M$, $\alpha = 95\%$).

| Method | <i>this paper</i> | <i>Rebaï</i> | <i>Feingold</i> |
|-----------|-------------------|--------------|-----------------|
| Threshold | 8.23 | 9.09 | 8.26 |

Fig. 5. Thresholds as a function of the method considered. The map consists of 51 genetic markers equally spaced every 2cM ($T = 1M$, $\alpha = 95\%$).

| Method | <i>this paper</i> | <i>Feingold</i> |
|-----------|-------------------|-----------------|
| Threshold | 5.40 | 5.78 |

Fig. 6. Thresholds as a function of the method considered. The map consists of 2 genetic markers ($T = 1M$, $\alpha = 95\%$).

5. Proofs

Notations : I_θ will be the Fisher information matrix taken at the point θ . $I_{ij}(\theta)$ refers to the element ij of I_θ . $I_{ij}^{-1}(\theta)$ refers to the element ij of I_θ^{-1} , the inverse of I_θ .

5.1. Proof of theorem 1

We first compute the score functions and the Fisher information matrix. Let $t \in [t_1, t_2]$.

$$\frac{\partial \log L}{\partial q} \Big|_{\theta_0} = \frac{y - \mu}{\sigma^2} \{2p(t) - 1\}$$

$$\frac{\partial \log L}{\partial \mu} \Big|_{\theta_0} = \frac{y - \mu}{\sigma^2} \quad , \quad \frac{\partial \log L}{\partial \sigma} \Big|_{\theta_0} = -\frac{1}{\sigma} + \frac{(y - \mu)^2}{\sigma^3}$$

$$I_{11}(\theta_0) = \frac{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}{\sigma^2} \quad , \quad I_{22}(\theta_0) = \frac{1}{\sigma^2}$$

As the fourth-order moment of a standard normal distribution is equal to three,

$$I_{33}(\theta_0) = \frac{2}{\sigma^2}$$

After some calculations, we find : $I_{12}(\theta_0) = I_{13}(\theta_0) = I_{23}(\theta_0) = 0$. So,

$$I_{\theta_0} = \text{Diag} \left[\frac{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}{\sigma^2} , \frac{1}{\sigma^2} , \frac{2}{\sigma^2} \right] \quad (8)$$

where $\mathbb{E} \left[\{2p(t_1) - 1\}^2 \right] = \mathbb{E} \left[\{2p(t_2) - 1\}^2 \right] = 1$ and $\forall t \in]t_1, t_2[:$

$$\mathbb{E} \left[\{2p(t) - 1\}^2 \right] = \bar{r}(t_1, t_2) \left(2Q_t^{1,1} - 1 \right)^2 + r(t_1, t_2) \left(2Q_t^{1,-1} - 1 \right)^2 \quad (9)$$

Indeed, $\forall t \in]t_1, t_2[:$

$$\begin{aligned} \mathbb{E} \left[\{2p(t) - 1\}^2 \right] &= 2 \left\{ \left(Q_t^{1,1} \right)^2 \bar{r}(t_1, t_2) + \left(Q_t^{1,-1} \right)^2 r(t_1, t_2) \right\} \\ &\quad + 2 \left\{ \left(Q_t^{-1,1} \right)^2 r(t_1, t_2) + \left(Q_t^{-1,-1} \right)^2 \bar{r}(t_1, t_2) \right\} - 1 \end{aligned}$$

As $Q_t^{-1,1} = 1 - Q_t^{1,-1}$, $Q_t^{-1,-1} = 1 - Q_t^{1,1}$ and $\bar{r}(t_1, t_2) + r(t_1, t_2) = 1$, we obtain formula (9).

$\mathbb{E} \left[\{2p(t) - 1\}^2 \right]$ is always different from zero since the parameter t is bounded. It comes $\forall t \in [t_1, t_2]$:

$$\Lambda_n(t) = \left[\sum_{j=1}^n \frac{(y_j - \mu) \{2p_j(t) - 1\}}{\sigma \sqrt{n} \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}} \right]^2 + o_{P_{\theta_0}}(1) \quad (10)$$

By convention, the notation $o_{P_{\theta_0}}(1)$ is short for a sequence of random vectors that converges to zero in probability under H_0 (i.e. no QTL on the whole interval studied).

Study under H_0 :

Without loss of generality, we assume that $n = 1$ for the moment and we consider the score function :

$$S(t) = \frac{(y - \mu) \{2p(t) - 1\}}{\sigma \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}} = \frac{y - \mu}{\sigma} h(t)$$

where the fact $h(\cdot)$ is a random process independent of y .

It is easy to see that :

$$\mathbb{E} \{S(t)\} = 0 \quad , \quad \mathbb{V} \{S(t)\} = \mathbb{E} \left[\{h(t)\}^2 \right] = 1$$

$\forall (t, t') \in [t_1, t_2]^2$:

$$\begin{aligned} \Gamma(t, t') := \text{Cov} \{S(t), S(t')\} &= \mathbb{E} \{h(t)h(t')\} = \frac{\mathbb{E} \left[\{2p(t) - 1\} \{2p(t') - 1\} \right]}{\sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]} \sqrt{\mathbb{E} \left[\{2p(t') - 1\}^2 \right]}} \\ &= \frac{4\mathbb{E} \{p(t)p(t')\} - 1}{\sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]} \sqrt{\mathbb{E} \left[\{2p(t') - 1\}^2 \right]}} \quad (11) \end{aligned}$$

The formula for $\mathbb{E} \{p(t)p(t')\}$ is given in appendix 7.1. As $|p(t)p(t')| \leq 1$, by dominated convergence theorem, $\mathbb{E} \{p(t)p(t')\}$ is continuous at (t_1, t') , (t_2, t') and (t_1, t_2) . Then the covariance function is continuous at this points (because the denominator is also continuous). So, the covariance function is a continuous function on $[t_1, t_2]^2$.

Let $S_n(\cdot)$ be the score process for n observations :

$$S_n(t) = \sum_{j=1}^n \frac{(y_j - \mu) (2p_j(t) - 1)}{\sigma \sqrt{n} \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}} \quad (12)$$

When n tends to infinity, an application of the Multivariate Central Limit Theorem shows that for $0 \leq s_1 < s_2 < \dots < s_d \leq T$:

$$(S_n(s_1), \dots, S_n(s_d))' \xrightarrow{\mathcal{L}} N(\mathbf{0}, \Sigma)$$

were Σ is the variance covariance matrix, with unit variance and covariance given by formula (11). $\underline{0}$ is a column vector of zeros. As $\Lambda_n(t) = S_n^2(t) + o_{P_{\theta_0}}(1)$:

$$(\Lambda_n(s_1), \dots, \Lambda_n(s_d))' \xrightarrow{\mathcal{L}} \left\{ N(\underline{0}, \Sigma) \right\}^2$$

Study under H_{at^*} :

In this part, we set

$$Y_j = \mu + \frac{a}{\sqrt{n}} X_j(t^*) + \sigma \varepsilon_j \quad (13)$$

where ε_j is a Gaussian white noise. According to formula (10), $\forall t \in [t_1, t_2]$:

$$\Lambda_n(t) = \{S_n(t)\}^2 + o_{P_{\theta_0}}(1)$$

We remind that $o_{P_{\theta_0}}(1)$ is short for a sequence of random vectors that converges to zero in probability under H_0 (i.e. no QTL on the whole interval studied). Let $o_{P_{\theta_0, t^*}}(1)$ be a sequence of random vectors that converges to zeros if there is no QTL at position t^* . Then, it is clear that :

$$\Lambda_n(t) = \{S_n(t)\}^2 + o_{P_{\theta_0, t^*}}(1)$$

Let θ_{a, t^*} be the parameter referring that we are under H_{at^*} . Under H_{at^*} , as the QTL is located at position t^* , the density of $Y|X(t_0), X(t_1)$ verifies :

$$p(t^*)f_{(\mu+q, \sigma)}(y) + \{1 - p(t^*)\}f_{(\mu-q, \sigma)}(y)$$

Let Q_n and P_n two sequences of probability measures defined on the same space $(\Omega_n, \mathcal{A}_n)$. Q_n (respectively P_n) is the law corresponding to the density $L_n(\theta_{a, t^*}, t^*)$ (resp $L_n(\theta_0, t^*)$).

We will call the log likelihood ratio $\log \frac{dQ_n}{dP_n}$. It verifies : $\log \frac{dQ_n}{dP_n} = \log \left\{ \frac{L_n(\theta_{a, t^*}, t^*)}{L_n(\theta_0, t^*)} \right\}$.

Notations : $Q_n \triangleleft P_n$ will mean the sequence Q_n is contiguous with the respect to the sequence P_n .

Let $b = (a, 0, 0)'$. As the model is differentiable in quadratic mean at θ_{a, t^*} :

$$\log \left(\frac{dQ_n}{dP_n} \right) = \frac{b'}{\sqrt{n}} \nabla \log L_n(\theta_0, t^*) - \frac{1}{2} b' I_{\theta_0} b + o_{P_{\theta_0}}(1)$$

Then, by the central limit theorem :

$$\log \left(\frac{dQ_n}{dP_n} \right) \xrightarrow{\mathcal{L}}_{H_0} N\left(-\frac{1}{2}\nu^2, \nu^2\right) \text{ with } \nu^2 = \frac{a^2}{\sigma^2} \mathbb{E} \left[\{2p(t^*) - 1\}^2 \right]$$

So, by the iii) of Le Cam's first lemma, we have $Q_n \triangleleft P_n$.

Up to now $\forall t \in [t_1, t_2]$:

$$\Lambda_n(t) = \{S_n(t)\}^2 + o_{P_{\theta_0, t^*}}(1)$$

As $Q_n \triangleleft P_n$, according to iv) of Le Cam's first lemma :

$$\Lambda_n(t) = \{S_n(t)\}^2 + o_{P_{\theta_{a,t^*}}}(1)$$

So, calculations can be done with the score process. According to formula (12) and (13), we have :

$$S_n(t) = \frac{1}{\sqrt{n}} \sum_{j=1}^n \varepsilon_j h_j(t) + \sum_{j=1}^n \frac{a}{\sigma n} X_j(t^*) h_j(t) = S_n^0(t) + \sum_{j=1}^n \frac{a}{\sigma n} X_j(t^*) h_j(t)$$

where $h_j(\cdot)$ is the equivalent of the process $h(\cdot)$ defined above but for the individual j . $S_n^0(\cdot)$ is the process obtained under H_0 .

By the law of large number :

$$\frac{1}{n} \sum_{j=1}^n X_j(t^*) h_j(t) \rightarrow \mathbb{E} \{X(t^*) h(t)\}$$

Let suppose $K = 2$ for the moment and, for example $(t, t^*) \in]t_1, t_2]^2$. Let us compute $\mathbb{E}[X(t^*) \{2p(t) - 1\}]$. We condition on $X(t_1)$ and $X(t_2)$. Consider, for example, the case $X(t_1) = X(t_2) = 1$. In this case, $p(t) = Q_t^{1,1}$ and we have :

$$\begin{aligned} \mathbb{E} [X(t^*) \{2p(t) - 1\} \mid X(t_1) = X(t_2) = 1] &= \mathbb{E} [X(t^*) \{2Q_t^{1,1} - 1\} \mid X(t_1) = X(t_2) = 1] \\ &= \{2Q_t^{1,1} - 1\} \mathbb{E} [X(t^*) \mid X(t_1) = X(t_2) = 1] \\ &= \{2Q_t^{1,1} - 1\} \left\{ \frac{\bar{r}(t_1, t^*) \bar{r}(t^*, t_2)}{\bar{r}(t_1, t_2)} - \frac{r(t_1, t^*) r(t^*, t_2)}{\bar{r}(t_1, t_2)} \right\} \\ &= \{2Q_t^{1,1} - 1\} \{Q_{t^*}^{1,1} - Q_{t^*}^{-1,-1}\} = \{2Q_t^{1,1} - 1\} \{2Q_{t^*}^{1,1} - 1\} \end{aligned}$$

Considering the four cases :

$$\begin{aligned} \mathbb{E} [X(t^*) \{2p(t) - 1\}] &= \{2Q_t^{1,1} - 1\} \{2Q_{t^*}^{1,1} - 1\} \frac{1}{2} \bar{r}(t_1, t_2) + \{2Q_t^{1,-1} - 1\} \{2Q_{t^*}^{1,-1} - 1\} \frac{1}{2} r(t_1, t_2) \\ &+ \{2Q_t^{-1,1} - 1\} \{2Q_{t^*}^{-1,1} - 1\} \frac{1}{2} r(t_1, t_2) + \{2Q_t^{-1,-1} - 1\} \{2Q_{t^*}^{-1,-1} - 1\} \frac{1}{2} \bar{r}(t_1, t_2) \\ &= \bar{r}(t_1, t_2) \{2Q_{t^*}^{1,1} - 1\} \{2Q_t^{1,1} - 1\} \\ &+ r(t_1, t_2) \{2Q_{t^*}^{1,-1} - 1\} \{2Q_t^{1,-1} - 1\} \end{aligned} \tag{14}$$

According to dominated convergence theorem, $\mathbb{E}[X(t^*) \{2p(t) - 1\}]$ is continuous on $[t_1, t_2]^2$. As a conclusion, $\forall (t, t^*) \in [t_1, t_2]^2$:

$$m_{t^*}(t) = \frac{a \mathbb{E}[X(t^*) \{2p(t) - 1\}]}{\sigma \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]}}$$

A non linear interpolation

After some calculations, we can remark that :

$$S_n(t) = \{ \alpha(t) S_n(t_1) + \beta(t) S_n(t_2) \} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

where $\text{Cov}_{H_0} \{S_n(t_1), S_n(t_2)\} = e^{-2t_2}$, $\alpha(t_1) = 1$, $\beta(t_1) = 0$, $\alpha(t_2) = 0$, $\beta(t_2) = 1$ and $\forall t \in]t_1, t_2[$:

$$\alpha(t) = Q_t^{1,1} + Q_t^{1,-1} - 1 \quad \text{and} \quad \beta(t) = Q_t^{1,1} - Q_t^{1,-1}$$

And it comes :

$$m_{t^*}(t) = \{ \alpha(t) m_{t^*}(t_1) + \beta(t) m_{t^*}(t_2) \} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

Weak convergence of the score process

To begin, we remind that $t_1 = 0$ and $t_2 = T$. As $p(t)$ and $\mathbb{E} \left[\{2p(t) - 1\}^2 \right]$ are continuous functions, each trajectory of the process $S_n(\cdot)$ is a continuous function on $[0, T]$. Let define the modulus of continuity of a continuous function x on $[0, T]$:

$$w_x(\delta) = \sup_{|t' - t| < \delta} |x(t') - x(t)| \quad \text{where} \quad 0 < \delta \leq T$$

According to theorem 8.2 of Billingsley (1999), the score process is tight if and only if the two following conditions hold :

- (a) the sequence $S_n(0)$ is tight.
- (b) For each positive ϵ and η , there exist a δ , with $0 < \delta < T$, and an integer n_0 such that $\mathbb{P} \{w_{S_n}(\delta) \geq \eta\} \leq \epsilon \quad \forall n \geq n_0$.

According to Prohorov, the sequence $S_n(0)$ is tight. So, a) is verified.

Let define the functions $\tilde{\alpha}(t)$ and $\tilde{\beta}(t)$ such as :

$$\tilde{\alpha}(t) = \alpha(t) / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}, \quad \tilde{\beta}(t) = \beta(t) / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

First, we can remark that $\forall \delta$ such as $0 < \delta \leq T$:

$$\begin{aligned} w_{S_n}(\delta) &= \sup_{|t' - t| < \delta} |S_n(t') - S_n(t)| \\ &= \sup_{|t' - t| < \delta} \left| \{ \tilde{\alpha}(t') - \tilde{\alpha}(t) \} S_n(t_1) + \{ \tilde{\beta}(t') - \tilde{\beta}(t) \} S_n(t_2) \right| \\ &\leq \max \{ |S_n(t_1)|, |S_n(t_2)| \} \left\{ w_{\tilde{\alpha}}(\delta) + w_{\tilde{\beta}}(\delta) \right\} \end{aligned}$$

Let $\epsilon > 0$ and $\eta > 0$, as the sequence $\max \{ |S_n(t_1)|, |S_n(t_2)| \}$ is uniformly tight, $\exists M$ such as $\forall n \geq 1 \quad \mathbb{P} \left[\max \{ |S_n(t_1)|, |S_n(t_2)| \} \geq M \right] \leq \epsilon$.

It comes, $\mathbb{P} \left[\max \{ |S_n(t_1)|, |S_n(t_2)| \} \left\{ w_{\tilde{\alpha}}(\delta) + w_{\tilde{\beta}}(\delta) \right\} \geq M \left\{ w_{\tilde{\alpha}}(\delta) + w_{\tilde{\beta}}(\delta) \right\} \right] \leq \epsilon$

As $\tilde{\alpha}(t)$ and $\tilde{\beta}(t)$ are continuous on the compact $[0, T]$, according to Heine's theorem,

these functions are uniformly continuous. So, let $v > 0$, $\exists \delta_1$ with $0 < \delta_1 < T$, such as $w_{\bar{\alpha}}(\delta_1) < v/2$ and $\exists \delta_2$ with $0 < \delta_2 < T$ such as $w_{\bar{\beta}}(\delta_2) < v/2$. Let $\delta = \min(\delta_1, \delta_2)$ then $w_{\bar{\alpha}}(\delta) + w_{\bar{\beta}}(\delta) < v$. If we impose $v = \eta/M$, then $\forall n \geq 1$, $\mathbb{P}\{w_{S_n}(\delta) \geq \eta\} \leq \epsilon$ which means b) of theorem 8.2 of Billingsley (1999) is fulfilled. So, the tightness of the score process is proved.

To conclude, the tightness and the convergence of finite-dimensional imply the weak convergence of the score process.

5.2. Proof of theorem 2

Introducing the process $\tilde{W}(\cdot)$:

We consider the process $\tilde{W}(\cdot)$ on $[0, 1]$ such as :

$$\tilde{W}(t) = \frac{(1-t)C_1 + tC_2}{\sqrt{(1-t)^2 + t^2 + 2\tilde{\rho}t(1-t)}}$$

We can remark that $\tilde{W}(0) = C_1$ and $\tilde{W}(1) = C_2$.

The interest is on the supremum of the process $\{\tilde{W}(\cdot)\}^2$:

$$\{\tilde{W}(t)\}^2 = \frac{(1-t)^2\{C_1\}^2 + 2t(1-t)C_1C_2 + t^2\{C_2\}^2}{(1-t)^2 + t^2 + 2\tilde{\rho}t(1-t)}$$

We will call respectively $N(t)$ and $D(t)$ the numerator and the denominator of the fraction above.

$$\frac{\partial N(t)}{\partial t} = 2\{(1-t)C_1 + tC_2\}\{C_2 - C_1\}$$

We can remark that :

$$(1-t)^2 + t^2 + 2t\tilde{\rho}(1-t) = 1 - 2(1-\tilde{\rho})t(1-t) \quad (15)$$

It comes : $\frac{\partial D(t)}{\partial t} = -2(1-\tilde{\rho})(1-2t)$. So,

$$\begin{aligned} \frac{\partial \{\tilde{W}(t)\}^2}{\partial t} &= [2\{(1-t)C_1 + tC_2\}\{C_2 - C_1\}\{1 - 2(1-\tilde{\rho})t(1-t)\} \\ &\quad + 2(1-\tilde{\rho})(1-2t)\{(1-t)C_1 + tC_2\}^2] / \{D(t)\}^2 \end{aligned}$$

We have :

$$\begin{aligned} \frac{\partial \{\tilde{W}(t)\}^2}{\partial t} &= 0 \\ \Leftrightarrow \{(1-t)C_1 + tC_2\} \\ &\times [\{C_2 - C_1\}\{1 - 2(1-\tilde{\rho})t(1-t)\} + (1-\tilde{\rho})(1-2t)\{(1-t)C_1 + tC_2\}] = 0 \end{aligned}$$

As $\{(1-t)C_1 + tC_2\}$ corresponds to a minimum, the focus is on the second term. This second term is equal to zero if :

$$\frac{C_2}{C_1} = \frac{1 + \tilde{\rho}}{1 + (\tilde{\rho} - 1)t} - 1$$

Let define the function $\psi_{\tilde{\rho}}(t)$ such as :

$$\psi_{\tilde{\rho}}(t) = \frac{1 + \tilde{\rho}}{1 + (\tilde{\rho} - 1)t} - 1$$

As $(\tilde{\rho} - 1)t$ is a decreasing function on $[0, 1]$, then $\psi_{\tilde{\rho}}(t)$ is an increasing function on $[0, 1]$ with $\psi_{\tilde{\rho}}(0) = \tilde{\rho}$ and $\psi_{\tilde{\rho}}(1) = \frac{1}{\tilde{\rho}}$.

Let define $\psi_{\tilde{\rho}}^{-1}$ the inverse function of $\psi_{\tilde{\rho}}$. After straightforward calculations, we find :

$$\psi_{\tilde{\rho}}^{-1}(u) = \frac{\tilde{\rho} - u}{(\tilde{\rho} - 1)(u + 1)}$$

So, the extremum between 0 and 1 is obtained for :

$$\tilde{\xi} = \frac{\tilde{\rho} C_1 - C_2}{(\tilde{\rho} - 1)\{C_2 + C_1\}}$$

After some calculations, using formula (15), we find that :

$$(1 - \tilde{\xi})^2 + \tilde{\xi}^2 + 2 \tilde{\xi} \tilde{\rho} (1 - \tilde{\xi}) = \frac{1 + \tilde{\rho}}{1 - \tilde{\rho}} \frac{\{C_1\}^2 + \{C_2\}^2 - 2 \tilde{\rho} C_1 C_2}{\{C_1 + C_2\}^2}$$

It comes :

$$\{\tilde{W}(\tilde{\xi})\}^2 = \frac{\{C_1\}^2 + \{C_2\}^2 - 2 \tilde{\rho} C_1 C_2}{(1 + \tilde{\rho})(1 - \tilde{\rho})}$$

So :

$$\begin{aligned} \sup_{t \in [0,1]} \{\tilde{W}(t)\}^2 &= \{\tilde{W}(\tilde{\xi})\}^2 \mathbf{1}_{\frac{C_2}{C_1} \in] \tilde{\rho}, \frac{1}{\tilde{\rho}} [} + \{C_2\}^2 \mathbf{1}_{\frac{C_2}{C_1} \in [\frac{1}{\tilde{\rho}}, +\infty [} \\ &+ \{C_1\}^2 \mathbf{1}_{\frac{C_2}{C_1} \in [0, \tilde{\rho}]} + \{C_1\}^2 \mathbf{1}_{\frac{C_2}{C_1} \in] -\infty, 0[\cap |C_1| > |C_2|} \\ &+ \{C_2\}^2 \mathbf{1}_{\frac{C_2}{C_1} \in] -\infty, 0[\cap |C_2| > |C_1|} \end{aligned} \quad (16)$$

A concise version of this formula is that the supremum of $\{\tilde{W}(\cdot)\}^2$ is the maximum of three random variables :

$$\sup_{t \in [0,1]} \{\tilde{W}(t)\}^2 = \max \left[\{C_1\}^2, \{\tilde{W}(\tilde{\xi})\}^2 \mathbf{1}_{\frac{C_2}{C_1} \in] \tilde{\rho}, \frac{1}{\tilde{\rho}} [}, \{C_2\}^2 \right]$$

Let $\gamma_1(t)$ and $\gamma_2(t)$ be two continuous functions on $[t_1, t_2]$. Besides, as in theorem 2, let suppose that $\frac{\gamma_2(t)}{\gamma_1(t) + \gamma_2(t)}$ is bounded by 0 and 1, and that these bounds are reached. We have $\forall t \in [t_1, t_2]$:

$$\tilde{W} \left(\frac{\gamma_2(t)}{\gamma_1(t) + \gamma_2(t)} \right) = D(t)$$

And,

$$\sup_{t \in [t_1, t_2]} \{D(t)\}^2 = \sup_{t \in [0,1]} \{\tilde{W}(t)\}^2$$

It concludes the proof.

5.3. Proof of lemma 1

Study of the supremum of the linear interpolated process $W(\cdot)$:

Let consider the process $W(\cdot)$ on $[t_1, t_2]$. It verifies $\forall t \in [t_1, t_2]$:

$$W(t) = \left\{ \frac{t_2 - t}{t_2 - t_1} W(t_1) + \frac{t - t_1}{t_2 - t_1} W(t_2) \right\} / \sqrt{\tau(t)}$$

where

$$\tau(t) = \left(\frac{t_2 - t}{t_2 - t_1} \right)^2 + 2 \frac{(t - t_1)(t_2 - t)}{(t_2 - t_1)^2} e^{-2(t_2 - t_1)} + \left(\frac{t - t_1}{t_2 - t_1} \right)^2$$

Using same notations as in theorem 2, let consider :

$$\gamma_1(t) = \frac{t_2 - t}{t_2 - t_1}, \quad \gamma_2(t) = \frac{t - t_1}{t_2 - t_1}, \quad C_1 = W(t_1), \quad C_2 = W(t_2), \quad \tilde{\rho} = e^{-2(t_2 - t_1)}$$

We will call $\gamma_3(t)$ the ratio $\frac{\gamma_2(t)}{\gamma_1(t) + \gamma_2(t)}$. We have $\gamma_3(t) = \frac{t - t_1}{t_2 - t_1}$. So, $\gamma_3(t_2) = 1$, $\gamma_3(t_1) = 0$, and $0 \leq \gamma_3(t) \leq 1$. As a consequence, according to theorem 2 :

$$\sup_{t \in [t_1, t_2]} \{W(t)\}^2 = \max \left[\{W(t_1)\}^2, \{W(t_2)\}^2, \frac{\{W(t_1)\}^2 + \{W(t_2)\}^2 - 2 e^{-2(t_2 - t_1)} W(t_1) W(t_2)}{\{1 + e^{-2(t_2 - t_1)}\} \{1 - e^{-2(t_2 - t_1)}\}} \mathbf{1}_{\frac{W(t_2)}{W(t_1)} \in] e^{-2(t_2 - t_1)}, e^{2(t_2 - t_1)} [} \right]$$

In the same way as in the proof of theorem 2 (cf. Section 5.2), we have :

$$\{W(\xi)\}^2 = \frac{\{W(t_1)\}^2 + \{W(t_2)\}^2 - 2 e^{-2(t_2 - t_1)} W(t_1) W(t_2)}{\{1 + e^{-2(t_2 - t_1)}\} \{1 - e^{-2(t_2 - t_1)}\}}$$

with :

$$\xi = \frac{(t_2 - t_1) \{e^{-2(t_2 - t_1)} W(t_1) - W(t_2)\}}{\{e^{-2(t_2 - t_1)} - 1\} \{W(t_1) + W(t_2)\}} + t_1$$

Study of the supremum of the non linear interpolated process $Z(\cdot)$:

Let consider the process $Z(\cdot)$ on $[t_1, t_2]$. It verifies $\forall t \in [t_1, t_2]$:

$$Z(t) = \{\alpha(t)Z(t_1) + \beta(t)Z(t_2)\} / \sqrt{\{\alpha(t)\}^2 + \{\beta(t)\}^2 + 2 \alpha(t) \beta(t) e^{-2(t_2 - t_1)}}$$

Indeed, according to the proof of theorem 1 (cf. Section 5.1), $Z(\cdot)$ has unit variance.

Using same notations as in theorem 2, let consider :

$$\gamma_1(t) = \alpha(t), \quad \gamma_2(t) = \beta(t), \quad C_1 = Z(t_1), \quad C_2 = Z(t_2), \quad \tilde{\rho} = e^{-2(t_2 - t_1)}$$

As previously, we will call $\gamma_3(t)$ the ratio $\frac{\gamma_2(t)}{\gamma_1(t) + \gamma_2(t)}$. We have $\gamma_3(t) = \frac{\beta(t)}{\alpha(t) + \beta(t)}$. To begin, we will admit that $\forall t \in [t_1, t_2]$, $0 \leq \gamma_3(t) \leq 1$. Besides, $\gamma_3(t_2) = 1$ and $\gamma_3(t_1) = 0$.

As a consequence, according to theorem 2 :

$$\sup_{t \in [t_1, t_2]} \{Z(t)\}^2 = \max \left[\{Z(t_1)\}^2, \{Z(t_2)\}^2, \frac{\{Z(t_1)\}^2 + \{Z(t_2)\}^2 - 2 e^{-2(t_2-t_1)} Z(t_1) Z(t_2)}{\{1 + e^{-2(t_2-t_1)}\} \{1 - e^{-2(t_2-t_1)}\}} 1^{\frac{Z(t_2)}{Z(t_1)} \in] e^{-2(t_2-t_1)}, e^{2(t_2-t_1)} [} \right]$$

To conclude, as $Z(t_1) = W(t_1)$ and $Z(t_2) = W(t_2)$:

$$\sup_{t \in [t_1, t_2]} \{Z(t)\}^2 = \sup_{t \in [t_1, t_2]} \{W(t)\}^2$$

The interest is now on the function $\gamma_3(t)$. $\forall t \in [t_1, t_2]$, we have :

$$\gamma_3(t) = 1 - \frac{Q_t^{1,1} + Q_t^{1,-1} - 1}{2Q_t^{1,1} - 1}$$

In order to prove that $\forall t \in]t_1, t_2[$ $0 < \gamma_3(t) < 1$, we will prove that $0 < \bar{\gamma}_3(t) < 1$ with $\bar{\gamma}_3(t) = (Q_t^{1,1} + Q_t^{1,-1} - 1)/(2Q_t^{1,1} - 1)$.

The calculation of the derivative of $Q_t^{1,1}$ shows that $Q_t^{1,1}$ is a decreasing function on $]t_1, (t_1 + t_2)/2[$ and an increasing function on $[(t_1 + t_2)/2, t_2[$. The minimum is $Q_{(t_1+t_2)/2}^{1,1}$. Since $Q_{(t_1+t_2)/2}^{1,1} = \{1 + e^{-2(t_1+t_2)} + 2e^{-t_1-t_2}\} / \{2 + 2e^{-2(t_2-t_1)}\}$, we have $Q_{(t_1+t_2)/2}^{1,1} > 1/2$. By continuity, $0 < 2Q_t^{1,1} - 1$. The focus is now on the numerator of $\bar{\gamma}_3(t)$. After calculations, we obtain :

$$\frac{\partial(Q_t^{1,1} + Q_t^{1,-1})}{\partial t} = \frac{-e^{-2(t-t_1)} - e^{-4t_2+2t+2t_1}}{2 r(t_1, t_2) \bar{r}(t_1, t_2)}$$

As a consequence, $Q_t^{1,1} + Q_t^{1,-1} - 1$ is a decreasing function on $]t_1, t_2[$. By continuity, $0 < Q_t^{1,1} + Q_t^{1,-1} - 1$. So, $\bar{\gamma}_3(t) > 0$.

On the other hand, the study of the derivative of $Q_t^{1,1} - Q_t^{1,-1}$ shows that $Q_t^{1,1} > Q_t^{1,-1}$. It comes $\forall t \in]t_1, t_2[$, $0 < \bar{\gamma}_3(t) < 1$ and $0 < \gamma_3(t) < 1$.

6. Acknowledgements

The authors thank Jean-Michel Elsen for having proposed this subject of research and fruitful discussions. This work has been supported by the Animal Genetic Department of the French National Institute for Agricultural Research, SABRE, and the National Center for Scientific Research.

7. Appendix

7.1. Formula for $\mathbb{E}\{p(t)p(t')\}$

$\forall (t, t') \in]t_1, t_2[^2$:

$$\begin{aligned} \mathbb{E}\{p(t)p(t')\} &= \frac{1}{2} \left\{ Q_t^{1,1} Q_{t'}^{1,1} \bar{r}(t_1, t_2) + Q_t^{1,-1} Q_{t'}^{1,-1} r(t_1, t_2) \right\} \\ &\quad + \frac{1}{2} \left\{ Q_t^{-1,1} Q_{t'}^{-1,1} r(t_1, t_2) + Q_t^{-1,-1} Q_{t'}^{-1,-1} \bar{r}(t_1, t_2) \right\} \end{aligned}$$

This quantity is continuous at t_1 and t_2 (cf. proof of theorem 1 in Section 5.1)

7.2. Sketch of the proof of theorem 3

Let $t \in [t_1, t_K] \setminus \mathbb{T}_k$. As t belongs to the "Marker interval" (t^ℓ, t^r) , some adjustments with Section 3 have to be done : t_1 becomes t^ℓ and t_2 becomes t^r . So, $p(t)$ is now the quantity equal to $\mathbb{P}\{X(t) = 1 | X(t^\ell), X(t^r)\}$. In the same way, $p(t)$, $Q_t^{1,1}$, $Q_t^{1,-1}$, $Q_t^{-1,1}$ and $Q_t^{-1,-1}$ described in formula (2) have to be adapted to the "Marker interval". The likelihood presented in formula (3), is unchanged except that the focus is on the triplet $(Y, X(t^\ell), X(t^r))$ and the function $g(t)$ has to be adapted to the "Marker interval". Formula (10) of Section 5.1 is also suitable $t \in [t_1, t_K] \setminus \mathbb{T}_k$ because t is bounded. It comes, $\forall (t, t') \in [t_1, t_K] \setminus \mathbb{T}_k \times [t_1, t_K] \setminus \mathbb{T}_k$:

$$\Gamma(t, t') = \frac{4\mathbb{E}\{p(t)p(t')\} - 1}{\sqrt{\mathbb{E}\{2p(t) - 1\}^2} \sqrt{\mathbb{E}\{2p(t') - 1\}^2}}$$

$\mathbb{E}\{2p(t) - 1\}^2$ described in formula (9) of Section 5.1 has to be adapted to the "Marker interval".

$\forall (t, t') \in]t^\ell, t^r[^2$, the expression of $\mathbb{E}\{p(t)p(t')\}$ can be deduced from appendix 7.1 by adapting to the "Marker interval".

Besides, if $(t, t') \in]t^\ell, t^r[\times [t^r, t_K] \setminus \mathbb{T}_k$:

$$\begin{aligned} & \mathbb{E}\{p(t)p(t')\} \\ &= \frac{1}{2}\bar{r}(t^\ell, t^r) \left[Q_{t'}^{1,1}\bar{r}\{(t')^\ell, (t')^r\} + Q_{t'}^{1,-1}r\{(t')^\ell, (t')^r\} \right] \left[Q_t^{1,1}\bar{r}\{t^r, (t')^\ell\} + Q_t^{-1,-1}r\{t^r, (t')^\ell\} \right] \\ &+ \frac{1}{2}\bar{r}(t^\ell, t^r) \left[Q_{t'}^{-1,1}r\{(t')^\ell, (t')^r\} + Q_{t'}^{-1,-1}\bar{r}\{(t')^\ell, (t')^r\} \right] \left[Q_t^{1,1}r\{t^r, (t')^\ell\} + Q_t^{-1,-1}\bar{r}\{t^r, (t')^\ell\} \right] \\ &+ \frac{1}{2}r(t^\ell, t^r) \left[Q_{t'}^{1,1}\bar{r}\{(t')^\ell, (t')^r\} + Q_{t'}^{1,-1}r\{(t')^\ell, (t')^r\} \right] \left[Q_t^{1,-1}r\{t^r, (t')^\ell\} + Q_t^{-1,1}\bar{r}\{t^r, (t')^\ell\} \right] \\ &+ \frac{1}{2}r(t^\ell, t^r) \left[Q_{t'}^{-1,1}r\{(t')^\ell, (t')^r\} + Q_{t'}^{-1,-1}\bar{r}\{(t')^\ell, (t')^r\} \right] \left[Q_t^{1,-1}\bar{r}\{t^r, (t')^\ell\} + Q_t^{-1,1}r\{t^r, (t')^\ell\} \right] \end{aligned}$$

In the same way as what has been done in the proof of theorem 1 (cf. Section 5.1),

$\forall (t, t^*) \in [t_1, t_K] \setminus \mathbb{T}_k \times [t_1, t_K] \setminus \mathbb{T}_k$:

$$m_{t^*}(t) = \frac{a \mathbb{E}[X(t^*)\{2p(t) - 1\}]}{\sigma \sqrt{\mathbb{E}\{2p(t) - 1\}^2}}$$

If $(t, t^*) \in]t^\ell, t^r[^2$, then $\mathbb{E}[X(t^*)\{2p(t) - 1\}]$ has the same expression as in formula (14) of Section 5.1 provided that we adapt to the "Marker interval".

Besides, if $(t, t^*) \in]t^\ell, t^r[\times [t^r, t_K] \setminus \mathbb{T}_k$:

$$\begin{aligned} & \mathbb{E}[X(t^*)\{2p(t) - 1\}] \\ &= 2 Q_t^{1,1} \mathbb{E}\{X(t^*)1_{X(t^\ell)=1}1_{X(t^r)=1}\} + 2 Q_t^{1,-1} \mathbb{E}\{X(t^*)1_{X(t^\ell)=1}1_{X(t^r)=-1}\} \\ &+ 2 Q_t^{-1,1} \mathbb{E}\{X(t^*)1_{X(t^\ell)=-1}1_{X(t^r)=1}\} + 2 Q_t^{-1,-1} \mathbb{E}\{X(t^*)1_{X(t^\ell)=-1}1_{X(t^r)=-1}\} \\ &= \bar{r}(t^\ell, t) \bar{r}(t, t^r) \{1 - 2r(t^r, t^*)\} + \bar{r}(t^\ell, t) r(t, t^r) \{2r(t^r, t^*) - 1\} \\ &+ r(t^\ell, t) \bar{r}(t, t^r) \{1 - 2r(t^r, t^*)\} + r(t^\ell, t) r(t, t^r) \{2r(t^r, t^*) - 1\} \\ &= \{1 - 2r(t, t^r)\} \{1 - 2r(t^r, t^*)\} = e^{-2(t^* - t)} \end{aligned}$$

As we deal with Poisson processes, it is reversible. So, If $(t, t^*) \in [t^{*r}, t_K] \setminus \mathbb{T}_k \times]t^{*\ell}, t^{*r}[$:

$$\mathbb{E}[X(t^*) \{2p(t) - 1\}] = \{1 - 2r(t^*, t^\ell)\} \{1 - 2r(t^\ell, t)\} = e^{-2(t-t^*)}$$

So, if t and t^* do not belong to the same "Marker interval" :

$$\mathbb{E}[X(t^*) \{2p(t) - 1\}] = e^{-2|t-t^*|} \quad (17)$$

A non linear interpolation

Concerning the non linear interpolation, we have to adapt formula (5) of Section 3.1 to the "Marker interval". $\forall t \in [t_1, t_K] \setminus \mathbb{T}_k$ we have :

$$S_n(t) = \{ \alpha(t) S_n(t^\ell) + \beta(t) S_n(t^r) \} / \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]} \quad (18)$$

where $\alpha(t) = Q_t^{1,1} + Q_t^{1,-1} - 1$, $\beta(t) = Q_t^{1,1} - Q_t^{1,-1}$ and $\forall k \forall k'$, $\text{Cov}_{H_0} \{S_n(t_k), S_n(t_{k'})\} = e^{-2|t_k - t_{k'}|}$.

It comes $\forall (t, t^*) \in [t_1, t_K] \setminus \mathbb{T}_k \times [t_1, t_K] \setminus \mathbb{T}_k$:

$$m_{t^*}(t) = \{ \alpha(t) m_{t^*}(t^\ell) + \beta(t) m_{t^*}(t^r) \} / \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]}$$

Weak convergence of the score process

Each trajectory of the process $S_n(\cdot)$ is a continuous function on $[0, T]$. In the same way as in the proof of theorem 1 in Section 5.1, in order to prove the tightness of the score process, we have to verify that conditions a) and b) of theorem 8.2 of Billingsley (1999) are fulfilled. According to Prohorov, $S_n(0)$ is tight, so a) is fulfilled.

We remind the modulus of continuity of $S_n(t)$:

$$w_{S_n}(\delta) = \sup_{|t'-t| < \delta} |S_n(t') - S_n(t)| \quad \text{where } 0 < \delta \leq T$$

Let define $w_{S_n}^k(\delta)$, the modulus of continuity of $S_n(t)$ only between the markers k and $k+1$:

$$w_{S_n}^k(\delta) = \sup_{|t'-t| < \delta} |S_n(t' + t_k) - S_n(t + t_k)| \quad \text{where } 0 < \delta \leq t_{k+1} - t_k$$

As the score process is tight when there are only two markers (cf. proof of theorem 1), according to b) of theorem 8.2 of Billingsley (1999), we have for a given k :

$$\forall \epsilon > 0 \forall \eta > 0 \exists \delta_k \text{ with } 0 < \delta_k < t_{k+1} - t_k \text{ such that } \mathbb{P}\{w_{S_n}^k(\delta_k) \geq \eta\} \leq \epsilon$$

So, let $\epsilon > 0$, $\epsilon' = \epsilon/(K-1)$, $\eta > 0$ and we impose $\delta = \min_{k \in \{1, \dots, K-1\}}(\delta_k)$

then $\forall k \in \{1, \dots, K-1\}$ $\mathbb{P}\{w_{S_n}^k(\delta) \geq \eta\} \leq \epsilon'$.

As $w_{S_n}(\delta) \geq w_{S_n}^1(\delta) + \dots + w_{S_n}^{K-1}(\delta)$, then $\mathbb{P}\{w_{S_n}(\delta) \geq \eta\} \leq \sum_{k=1}^{K-1} \mathbb{P}\{w_{S_n}^k(\delta) \geq \eta\} \leq \epsilon$ which means b) of theorem 8.2 of Billingsley (1999) is fulfilled. So, the tightness of the score process is proved.

To conclude, the tightness and the convergence of finite-dimensional imply the weak convergence of the score process.

7.3. Linear interpolated process in presence of several markers

In presence of several markers, the process $V_n(\cdot)$ is such as $\forall t \in [t_1, t_K] \setminus \mathbb{T}_k$:

$$V_n(t) = \left\{ \frac{t^r - t}{t^r - t^\ell} S_n(t^\ell) + \frac{t - t^\ell}{t^r - t^\ell} S_n(t^r) \right\} / \sqrt{\tau(t)}$$

where

$$\tau(t) = \left(\frac{t^r - t}{t^r - t^\ell} \right)^2 + 2 \frac{(t^r - t)(t - t^\ell)}{(t^r - t^\ell)^2} e^{-2(t^r - t^\ell)} + \left(\frac{t - t^\ell}{t^r - t^\ell} \right)^2$$

It can be seen easily that $\tau(t) \neq 0$, $\forall t \in [t_1, t_K] \setminus \mathbb{T}_k$.

$V_n(\cdot)$ remains asymptotically a Gaussian process with mean equal to 0 under H_0 , unit variance, and $\forall k \forall k'$, $\text{Cov}_{H_0} \{S_n(t_k), S_n(t_{k'})\} = e^{-2|t_k - t_{k'}|}$. In the same way as what has been done in Section 3.3, the weights of the model of mixture corresponding to this process verify :

$$p(t) = 1_{X(t^\ell)=1} 1_{X(t^r)=1} + \frac{t^r - t}{t^r - t^\ell} 1_{X(t^\ell)=1} 1_{X(t^r)=-1} + \frac{t - t^\ell}{t^r - t^\ell} 1_{X(t^\ell)=-1} 1_{X(t^r)=1}$$

This weights are an approximation at the first order of the original weights. So, $V_n(\cdot)$ will be a good approximation if and only if the genetic markers are close to each other. This process $V_n(\cdot)$ is a generalization of the process studied, under H_0 , by Rebaï et al. (1994). By contiguity (in the same way of what has been done in Section 5.1), under H_{at^*} , $V_n(\cdot)$ is asymptotically the same process as under H_0 on which the mean function, $\tilde{m}_{t^*}(t)$, has been added :

$$\tilde{m}_{t^*}(t) = \left\{ \frac{t^r - t}{t^r - t^\ell} m_{t^*}(t^\ell) + \frac{t - t^\ell}{t^r - t^\ell} m_{t^*}(t^r) \right\} / \sqrt{\tau(t)}$$

As previously, $W(\cdot)$, the limiting process of $V_n(\cdot)$, is named the linear interpolated process.

7.4. Comparison with Chang et al. (2009)

The law of the LRT process has also been obtained by Chang et al. (2009) under the null hypothesis. We propose here to present technical differences between our work and the work of Chang et al. (2009). As at a location t , the LRT is asymptotically the square of the score test, we will focus only on the score process as in Chang et al. (2009).

The main difference between the two approaches is that we consider the number of individuals in each class as a random variable whereas in Chang et al. (2009), the number of individuals in each class is supposed equal to the expectations (same remark as (b) of Section 3.3).

Our approach allows us to compute the score function $\frac{\partial \log L}{\partial q} |_{\theta_0}$ for only one observation and to calculate the Fisher information matrix without approximation.

Anyway, we obtain exactly the same Fisher information matrix as in Chang et al. (2009). However, there are some differences concerning other quantities.

7.4.1. Only two markers :

Let consider that there is only two markers as described in Section 3. Let $t \in]t_1, t_2[$. The result will be prolonged by continuity at the markers positions. According to formula (4)

of Section 3.3 and using the fact that $Q_t^{1,1} = 1 - Q_t^{-1,-1}$ and $Q_t^{1,-1} = 1 - Q_t^{-1,1}$, the score test statistic is :

$$S_n(t) = (1 - 2Q_t^{-1,-1}) \sum_{j=1}^n \frac{(y_j - \mu) \{1_{X_j(t_1)=1} 1_{X_j(t_2)=1} - 1_{X_j(t_1)=-1} 1_{X_j(t_2)=-1}\}}{\sigma \sqrt{n} \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]}} \\ + (1 - 2Q_t^{-1,1}) \sum_{j=1}^n \frac{(y_j - \mu) \{1_{X_j(t_1)=1} 1_{X_j(t_2)=-1} - 1_{X_j(t_1)=-1} 1_{X_j(t_2)=1}\}}{\sigma \sqrt{n} \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]}}$$

With our notations, the test statistic used in formula (8) of Chang et al. (2009) is :

$$U^*(t) = \frac{\sqrt{n}}{2} (1 - 2Q_t^{-1,-1}) \frac{\bar{r}(t_1, t_2) (\bar{y}_{11} - \bar{y}_{-1-1})}{\hat{\sigma} \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]}} + \frac{\sqrt{n}}{2} (1 - 2Q_t^{-1,1}) \frac{r(t_1, t_2) (\bar{y}_{1-1} - \bar{y}_{-11})}{\hat{\sigma} \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]}}$$

where $\bar{y}_{11} = \frac{2}{n\bar{r}(t_1, t_2)} \sum_{j=1}^n 1_{X_j(t_1)=1} 1_{X_j(t_2)=1}$, $\bar{y}_{-1-1} = \frac{2}{nr(t_1, t_2)} \sum_{j=1}^n 1_{X_j(t_1)=-1} 1_{X_j(t_2)=-1}$
 $\bar{y}_{1-1} = \frac{2}{nr(t_1, t_2)} \sum_{j=1}^n 1_{X_j(t_1)=1} 1_{X_j(t_2)=-1}$ and $\bar{y}_{-11} = \frac{2}{n\bar{r}(t_1, t_2)} \sum_{j=1}^n 1_{X_j(t_1)=-1} 1_{X_j(t_2)=1}$.

We can remark $S_n(t) \neq U^*(t) + o_{P_{\theta_0}}(1)$. It is due to the approximations done by Chang et al. (2009).

Let $G_n^1(t)$ and $G_n^2(t)$ be the quantities such as :

$$G_n^1(t) = \sum_{j=1}^n \frac{(y_j - \mu) \{1_{X_j(t_1)=1} 1_{X_j(t_2)=1} - 1_{X_j(t_1)=-1} 1_{X_j(t_2)=-1}\}}{\sigma \sqrt{n} \bar{r}(t_1, t_2)} \\ G_n^2(t) = \sum_{j=1}^n \frac{(y_j - \mu) \{1_{X_j(t_1)=1} 1_{X_j(t_2)=-1} - 1_{X_j(t_1)=-1} 1_{X_j(t_2)=1}\}}{\sigma \sqrt{n} r(t_1, t_2)}$$

$G_n^1(t)$ and $G_n^2(t)$ are asymptotically standard normal variables under H_0 . Besides, $G_n^1(t)$ and $G_n^2(t)$ are independent. Note that $G_n^1(t)$ and $G_n^2(t)$ do not depend on t but we keep t as a parameter in order to adapt these test statistics to the case of several markers in the next Section.

Contrary to formula (9) of Chang et al. (2009) :

$$G_n^1(t) \neq \frac{1}{2} \sqrt{\bar{r}(t_1, t_2)n} \frac{\bar{y}_{11} - \bar{y}_{-1-1}}{\hat{\sigma}} + o_{P_{\theta_0}}(1) \\ G_n^2(t) \neq \frac{1}{2} \sqrt{r(t_1, t_2)n} \frac{\bar{y}_{1-1} - \bar{y}_{-11}}{\hat{\sigma}} + o_{P_{\theta_0}}(1)$$

We have :

$$S_n(t) = \left\{ \sqrt{\bar{r}(t_1, t_2)} (1 - 2Q_t^{-1,-1}) G_n^1(t) + \sqrt{r(t_1, t_2)} (1 - 2Q_t^{-1,1}) G_n^2(t) \right\} / \sqrt{\mathbb{E}[\{2p(t) - 1\}^2]} \quad (19)$$

This formula is the corrected version of formula (10) of Chang et al. (2009) without approximations here. According to formula (19), the score at a position t between two markers,

is an interpolation not linear between the test statistic $G_n^1(t)$ and $G_n^2(t)$. Naturally, when t tends to t_1 (resp. t_2), $S_n(t)$ tends to $S_n(t_1)$ (resp. $S_n(t_2)$). It becomes a linear interpolation between $S_n(t_1)$ and $S_n(t_2)$ if a Taylor linearization is done concerning the weights of the model of mixture (cf. Section 3.3).

Finally, we agree with formula (11) of Chang et al. (2009) concerning the covariance of the process, it is exactly the same function as $\Gamma(t, t')$ of theorem 1 of this paper.

Note that the non linear interpolation presented above, in formula (19), is not the same interpolation as presented in formula (5) of Section 3.2 of this paper. Our interpolation is more intuitive, because it is an interpolation between the test statistic on markers. Besides, it explains why the likelihood profiles (ie. the paths of the process $\Lambda_n(\cdot)$) are smooth between markers.

7.4.2. Several markers : the ‘‘Interval Mapping’’ of Lander and Botstein (1989)

Let consider that there are several markers as described in Section 4. We consider values t, t' of the parameters that are distinct of markers positions. Let $t \in [t_1, t_K] \setminus \mathbb{T}_k$. We have :

$$G_n^1(t) = \sum_{j=1}^n \frac{(y_j - \mu) \{1_{X_j(t^\ell)=1} 1_{X_j(t^r)=1} - 1_{X_j(t^\ell)=-1} 1_{X_j(t^r)=-1}\}}{\sigma \sqrt{n \bar{r}(t^\ell, t^r)}}$$

$$G_n^2(t) = \sum_{j=1}^n \frac{(y_j - \mu) \{1_{X_j(t^\ell)=1} 1_{X_j(t^r)=-1} - 1_{X_j(t^\ell)=-1} 1_{X_j(t^r)=1}\}}{\sigma \sqrt{n r(t^\ell, t^r)}}$$

$$S_n(t) = \left\{ \sqrt{\bar{r}(t^\ell, t^r)} (2Q_t^{1,1} - 1) G_n^1(t) + \sqrt{r(t^\ell, t^r)} (2Q_t^{1,-1} - 1) G_n^2(t) \right\} / \sqrt{\mathbb{E} \left[\{2p(t) - 1\}^2 \right]}$$

This last formula is the corrected version of formula (14) of Chang et al. (2009).

Let $(t, t') \in]t^\ell, t^r[\times [t^r, t_K] \setminus \mathbb{T}_k$. The different covariances under H_0 are :

$$\begin{aligned} \text{Cov}_{H_0} \{G_n^1(t), G_n^1(t')\} &= \sqrt{\bar{r}(t^\ell, t^r) \bar{r} \{(t')^\ell, (t')^r\}} e^{-2\{(t')^\ell - t^r\}} \\ \text{Cov}_{H_0} \{G_n^1(t), G_n^2(t')\} &= \sqrt{\bar{r}(t^\ell, t^r) r \{(t')^\ell, (t')^r\}} e^{-2\{(t')^\ell - t^r\}} \\ \text{Cov}_{H_0} \{G_n^2(t), G_n^1(t')\} &= -\sqrt{r(t^\ell, t^r) \bar{r} \{(t')^\ell, (t')^r\}} e^{-2\{(t')^\ell - t^r\}} \\ \text{Cov}_{H_0} \{G_n^2(t), G_n^2(t')\} &= -\sqrt{r(t^\ell, t^r) r \{(t')^\ell, (t')^r\}} e^{-2\{(t')^\ell - t^r\}} \end{aligned}$$

This is exactly the same covariances as in formula (19) of Chang et al. (2009). Besides, we agree with formula (20) of Chang et al. (2009) which establish a relationship between the test statistic G when t and t' belong to 2 consecutive marker interval (as above we suppose $t < t'$):

$$G_n^2(t') = \frac{1}{\sqrt{r(t^r, (t')^r)}} \left\{ \sqrt{\bar{r}(t^\ell, t^r)} G_n^1(t) - \sqrt{r(t^\ell, t^r)} G_n^2(t) - \sqrt{\bar{r}(t^r, (t')^r)} G_n^1(t') \right\}$$

To conclude, the non linear interpolation proposed by Chang et al. (2009) is an approximation. We present here their interpolation without approximations. However, their approximations don't affect the final results concerning the process.

References

- Azaïs, J. M. and Cierco-Ayrolles, C. (2002). An asymptotic test for quantitative gene detection. *Ann. I. H. Poincaré*, **38**, **6**, 1087-1092.
- Azaïs, J. M., Gassiat, E., Mercadier, C. (2006). Asymptotic distribution and local power of the likelihood ratio test for mixtures. *Bernoulli*, **12**(5), 775-799.
- Azaïs, J. M., Gassiat, E., Mercadier, C. (2009). The likelihood ratio test for general mixture models with possibly structural parameter. *ESAIM*, To appear.
- Azaïs, J. M. and Wschebor, M. (2009). *Level sets and extrema of random processes and fields*. Wiley, New-York.
- Billingsley, P. (1999). *Convergence of probability measures*. Wiley, New-York.
- Chang, M. N., Wu, R., Wu, S. S., Casella, G. (2009). Score statistics for mapping quantitative trait loci. *Statistical Application in Genetics and Molecular Biology*, **8**(1), 16.
- Cierco, C. (1998). Asymptotic distribution of the maximum likelihood ratio test for gene detection. *Statistics*, **31**, 261-285.
- Davies, R.B. (1977). Hypothesis testing when a 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.
- Feingold, E., Brown, P.O., Siegmund, D. (1993). Gaussian models for genetic linkage analysis using complete high-resolution maps of identity by descent. *Am. J. Human. Genet.*, **53**, 234-251.
- Genz, A. (1992). Numerical computation of multivariate normal probabilities. *J. Comp. Graph. Stat.*, 141-149.
- Haldane, J.B.S (1919). The combination of linkage values and the calculation of distance between the loci of linked factors. *Journal of Genetics*, **8**, 299-309.
- Lander, E.S., Botstein, D. (1989). Mapping mendelian factors underlying quantitative traits using RFLP linkage maps. *Genetics*, **138**, 235-240.
- Le Cam, L. (1986). *Asymptotic Methods in Statistical Decision Theory*, Springer.
- Rabier, C-E. (2010). *PhD thesis*, Université Toulouse 3, Paul Sabatier.
- Rebaï, A., Goffinet, B., Mangin, B. (1994). Approximate thresholds of interval mapping tests for QTL detection. *Genetics*, **138**, 235-240.
- Rebaï, A., Goffinet, B., Mangin, B. (1995). Comparing power of different methods for QTL detection. *Biometrics*, **51**, 87-99.
- Siegmund, D. (1985). Sequential analysis : tests and confidence intervals. *Springer, New York*.

Van der Vaart, A.W. (1998) *Asymptotic statistics*, Cambridge Series in Statistical and Probabilistic Mathematics.

Wu, R., MA, C.X., Casella, G. (2007) *Statistical Genetics of Quantitative Traits*, Springer