

---

**Inferring genomic loss and location of tumor suppressor genes from high density genotypes**

Hui Wang<sup>1</sup>, Yohan Lee<sup>2</sup>, Stanley Nelson<sup>2</sup>, and Chiara Sabatti<sup>1,2</sup>

Departments of Statistics<sup>1</sup> and Human Genetics<sup>2</sup>, UCLA, Los Angeles CA 90095-7088

---

UCLA Statistics Department Preprint # 423

April 2005



**Running head** LOH studies with genotyping arrays

**Keywords** Loss of heterozygosity; single nucleotide polymorphisms; genotyping arrays.

**Corresponding author** Chiara Sabatti

Department of Human Genetics

UCLA School of Medicine

695 Charles E. Young Drive South

Los Angeles, California 90095-7088 (USA)

FAX: (310) 794-5446

Phone: (310) 794-9567

e-mail: [csabatti@mednet.ucla.edu](mailto:csabatti@mednet.ucla.edu)

## Abstract

Novel technologies, such as the 10k Affymetrix genotyping array, allow scoring of genetic polymorphisms at a very high density across the genome. This allows researchers to conduct traditional inquiries at an unprecedented resolution, while simultaneously motivates novel types of analysis, aimed at exploiting the increased information contained in these datasets. We consider how genotypes of cancer cell lines can be used to reconstruct genomic loss events and map putative tumor suppressor genes (TSG). Using a hidden Markov model framework, we adapt a previously described model for genomic instability in cancers to the current data structure. Simulations indicate that our procedure can be powerful and accurate and initial application to real data leads to encouraging results.

## 1 Introduction

Large scale genomic variation is receiving increasing attention from the scientific community (Pollack *et al.*, 2002; Bignell *et al.*, 2004; Cox *et al.*, 2005; Iafrate *et al.*, 2004; Sebat *et al.*, 2004). Along the process that leads to tumorigenesis, genomic stability is impaired: cancer cells present a higher rate of genomic losses and duplications. Moreover, some of these variations in copy number may determine the cancerous status of the cells, for instance by inactivating a tumor suppressor gene. Observed cancer cells are often considered to be selected to favor genomic losses in particular regions harboring tumor suppressor genes (TSG). Indeed, genomic losses have often been studied in cancer cells with the precise intent of locating TSGs (Newton *et al.* (1998) provides a comprehensive review of the biological process). Generally, a link between copy number variation and altered behaviour of cancer cells has been established (Pollack *et al.* 2002, Cox *et al.*, 2005). Concurrently, it has recently become apparent that there are a number of genetic polymorphisms that consist of large scale genomic variations (e.g. considerably sized genomic deletes and duplicates) that arise or are transmitted through the germline (Iafrate *et al.*, 2004; Sebat *et al.*, 2004).

While these two biological processes are very different, their effects can be measured with the same technology. For a number of years, large scale studies of genomic instability were limited to the case of cancers, and variation in copy numbers was assessed by genotyping tumor tissue and a normal cell from the same individual. In these studies, the fact that a marker is heterozygous in the normal cell and a homozygous in the cancer cell, can be considered as evidence that one of the chromosomes of the individual under study experienced a loss of genetic material and, hence, only the allele residing on the other chromosome remains to be detected. (Symmetrically, one could explain this result in terms of increase in copy number of one of the chromosomes, such that the signal from the duplicated allele overwhelms the one from the other). Newton *et al.* (1998) (Newton and Lee, 2000; Newton, 2002, Newton *et al.*, 2003) propose a model for the analysis of this loss of heterozygosity (LOH) data and demonstrate how to apply it to locate putative tumor suppressor genes.

Recently, a number of other technologies have been developed that enable a more direct and high throughput assessment of copy number variations: these include comparative genomic hybridization (CGH), and array-based CGH (Pinkel *et al.*, 1998; Bignell *et al.* 2004). Due to these technological advancements, scientists are able to observe recently documented large scale genomic variation (Iafate *et al.*, 2004; Sebat *et al.*, 2004). A number of statistical methods have been proposed to analyze the data from these experiments (Fridlyand *et al.*, 2004; Lai and Zhao, 2005; Wang *et al.*, 2005). Simultaneously, genotyping technology has evolved, and we are now able to measure polymorphic sites at a much higher density than just a few years ago. Thus LOH remains a powerful method of investigation. For example, recent studies have shown how genotyping arrays can be used to conduct effective copy number investigations (Lim *et al.*, 2004). One advantage of high density genotyping is that one can potentially gather information on copy number without the concurrent need to type both a “normal” and “case” cell from the same individual. This is due to the abundance of markers which provide a fine resolution snapshot of the genome, in which a “long” contiguous stretch of homozygous SNPs calls can be interpreted as the result of a

genomic alteration (with appropriate caveats). When one considers this potential (already noted in Lim *et al.*, 2004), developing methods that enable the use of high density genotypes to study copy number variations becomes an important goal. While other techniques may lead to more direct and precise measurements, and should likely be preferred if the primary goal of the study is detection of copy number variations, genotyping is done routinely on large scale both for tumor and normal cells. Failing to identify detectable changes in copy numbers from these data may result in erroneous interpretation of the results of the study, in addition to a loss of useful information.

In this paper, we focus on the use of high density SNP genotyping for the study of genomic instability in cancer cells. In a companion study (Wang *et al.*, 2005) we analyze the case of large scale genomic variation. In the next section we present a model for genotypes in the presence of genomic instability and selection. We subsequently illustrate how to estimate the instability parameters of the model and how to reconstruct the most likely profiles of genomic alteration from an individual's genotype data. Section 4 contains the description of a likelihood ratio test to identify location of a tumor suppressor gene, and Section 5 illustrates our results with simulated and real data. We conclude with a discussion.

## **2 A model for genotypes under genomic instability and selection**

We assume that the results of high density SNP genotyping across the genome are available for  $T$  cancer cell lines. In particular, we consider the case where the 10k genotyping array from Affymetrix is used, leading to the scoring of 10913 markers across the genome, at approximately intervals of 0.3 megabases. The methodology we will describe can be used with allele calls produced with other platforms, but it is tailored for high density genotyping—necessary to inform the inference of genomic aberrations from the genotype of cancer cells only. We denote with

$Y = \{y_i\}_{i=1}^M$  the sequence of genotypes at  $M$  markers for one cell line: each  $y_i$  can take on one of four possible values:  $(AA, AB, BB, -)$ , corresponding to the three genotypes and a “no call” value. We will assume independence between genotypes corresponding to markers on different chromosomes. It is convenient to group markers according to the chromosome where they are located:  $X_k = \{x_{ki}\}_{i=1}^{m_k}$  will indicate the genotypes of the  $m_k$  markers on chromosome  $k$ , so that  $Y = \{X_k\}_{k=1}^22$ . When unnecessary for clarity, we will avoid using the index  $k$ , so that  $X = \{x_i\}_{i=1}^m$  the collection of genotypes for markers in an unspecified chromosome. We use the superscript  $t$  to identify cell lines, so that the entire collection of our data is  $\{Y^t\}_{t=1}^T$ .

Newton *et al.* (1998) and Newton and Lee (2000) describe the model for genomic instability and selection in cancer cells that we use for our analysis. We refer the interested reader to these original papers to fully appreciate the underlying biological hypotheses, and we limit ourselves to a brief exposition. The substantial difference between the work presented in Newton *et al.* (1998), Newton and Lee (2000) and ours study resides in the fact that these previous authors asserted LOH calls, from the availability of genotypes from both a cancer and a normal cell from the same individual. Here we assume that only cancer cells are typed: the genomic loss process is unobserved, and we need to describe how it is reflected in the genotypes of cancer cells. The framework of hidden Markov models appears as a natural solution. We start considering the instability component of the model, and data from one cell line. Let  $\Pi = \{\pi_i\}_{i=1}^m$  denote the genomic loss process at the positions corresponding to each of the  $m$  markers on a chromosome:  $\pi_i = 0$  indicates no genomic alterations at the location of the  $i$ -th marker and  $\pi_i = 1$  indicates an alteration. In the following we will refer to alterations as a loss, since the model by Newton *et al.* (1998) that we are adapting was conceived specifically for losses; however, it is quite possible to interpret the detected abnormalities as increases in genomic copy number. With  $d_i$  we indicate the distance in megabases (Mb) between marker  $i$  and  $i + 1$ . Thus set the notation, the transition probabilities of

the hidden Markov process can be described in terms of two parameters  $\delta$  and  $\eta$ :

$$\begin{pmatrix} t(\pi_{i+1} = 1 \mid \pi_i = 1) & t(\pi_{i+1} = 0 \mid \pi_i = 1) \\ t(\pi_{i+1} = 1 \mid \pi_i = 0) & t(\pi_{i+1} = 0 \mid \pi_i = 0) \end{pmatrix} = \begin{pmatrix} 1 - (1 - \delta)(1 - e^{-\eta d_i}) & (1 - \delta)(1 - e^{-\eta d_i}) \\ \delta(1 - e^{-\eta d_i}) & 1 - \delta(1 - e^{-\eta d_i}) \end{pmatrix}.$$

The parameter  $\delta$  represents the sporadic loss rate, that is the probability that any location in the genome is lost in a random individual. The parameter  $\eta$  is used to model the dependency of the Markov process and the length of the genomic losses. In this framework the distance between two change-points in the  $\pi$  process is modeled similarly to the distance between two recombination events (Lange, 2002). Note that  $t(\pi_{i+1} = 1 \mid \pi_i = 1) \rightarrow 1$  and  $t(\pi_{i+1} = 0 \mid \pi_i = 0) \rightarrow 1$  as  $d_i \rightarrow 0$ . The adequacy of this model to describe genomic instability in cancer cells has been discussed by Newton *et al.* (1998) and supported by its successful application in empirical studies (see, for example Miller *et al.*, 2003). In this model,  $\delta$  and  $\eta$  are constant across a region spanned by linked markers. Depending on the nature of the data acquired, it may be sensible to assume that each chromosome is characterized by a specific value of  $\delta$  and  $\eta$ ; additionally, it may be appropriate to allow these parameters to be cell-line specific.

To link the unobserved loss process to the genotype data, we use the following emission probabilities:

$$\begin{pmatrix} e(x_i = AA \mid \pi_i = 1) & e(x_i = AB \mid \pi_i = 1) & e(x_i = BB \mid \pi_i = 1) & e(x_i = - \mid \pi_i = 1) \\ e(x_i = AA \mid \pi_i = 0) & e(x_i = AB \mid \pi_i = 0) & e(x_i = BB \mid \pi_i = 0) & e(x_i = - \mid \pi_i = 0) \end{pmatrix} = \begin{pmatrix} P_{A,i}(1 - \tau) & 0 & (1 - P_{A,i})(1 - \tau) & \tau \\ P_{A,i}^2(1 - \kappa) & 2P_{A,i}(1 - P_{A,i})(1 - \kappa) & (1 - P_{A,i})^2(1 - \kappa) & \kappa \end{pmatrix}, \quad (1)$$

where  $P_{A,i}$  is the frequency of allele  $A$  for the  $i$ th marker,  $\tau$  is the missing rate in loss regions, and  $\kappa$  is the missing rate in regions with no genomic aberrations. The difference in missing rate is due to the fact that the loss region may produce an ‘‘anomalous’’ intensity signal, filtered out by quality control mechanisms and resulting in ‘‘no call’’s. In the presence of a genomic loss, only homozygous genotypes are observed, and the relative abundance of  $AA$  and  $BB$  depends on the

allele frequencies  $P_{A,i}$ . While the emission probabilities (1) do not account for genotyping error, this effect can be easily incorporated. More substantial modification of the emission probabilities are needed, instead, to consider linkage disequilibrium across markers. We discuss such modifications in a companion paper (Wang *et al.*, 2005). The likelihood of a genotype sequence  $X$  under the instability model can be evaluated with standard HMM recursion formulas. In particular, if we define  $\alpha(\pi_i) = \Pr(x_1, \dots, x_i, \pi_i)$ , and  $\beta(\pi_i) = \Pr(x_{i+1}, \dots, x_m | \pi_i)$ , we have:

$$\begin{aligned}\alpha(\pi_i) &= \sum_{\pi_{i-1}=0,1} \alpha(\pi_{i-1})t(\pi_i|\pi_{i-1})e(x_i|\pi_i) \\ \beta(\pi_i) &= \sum_{\pi_{i+1}=0,1} \beta(\pi_{i+1})t(\pi_{i+1}|\pi_i)e(x_{i+1}|\pi_{i+1}).\end{aligned}$$

Then, for example  $\Pr(X) = \sum_{\pi_m=0,1} \alpha(\pi_m)$ . To emphasize that this probability depends only on the instability component of the model, we will indicate it with  $P_I(X)$ . A version of the aforementioned recursion formulas can also be used to evaluate conditional probabilities  $P_I(X|\pi_s = k)$ , which will be relevant in the following.

Newton *et al.* (1998) describe the selection effect that is the basis of the possible localization of tumor suppressor genes, and we adopt the same model. They consider the possibility of one tumor suppressor gene per chromosome. Two parameters are introduced:  $s$  represents the location of the tumor suppressor gene, and  $\omega_s$  the probability that a cancer cell line has a loss at position  $s$ . The likelihood of a genotype sequence, once the selection effect is introduced, can be written as:

$$\Pr(X) = \omega_s P_I(X|\pi_s = 1) + (1 - \omega_s) P_I(X|\pi_s = 0),$$

where  $P_I(X|\pi_s)$  depends on the instability process. In order to map a tumor suppressor gene, one needs to acquire data on multiple cell lines, so that the complete data for analysis will be a collection of genotype sequences  $X^1, X^2, \dots, X^T$  on the same chromosome. The data likelihood is then:

$$\mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega, \tau, s) = \prod_{t=1}^T (\omega P_I(X^t | \pi_s^t = 1) + (1 - \omega) P_I(X^t | \pi_s^t = 0)).$$

Given  $\mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega, \tau, s)$ , we can estimate the model parameters, reconstruct the location of most likely losses, and attempt mapping the tumor suppressor gene. We will assume that allele frequencies  $P_{A,i}$  are known: Affymetrix, for example, provides an estimate of allele frequencies for all the SNPs on its genotyping chip. We will also assume a known background “NoCall” rate: typically this can be estimated in each laboratory without much difficulty using a set of genotypes from normal individuals. The parameters we then need to estimate are  $\eta, \delta, \tau$  for the instability component of our model and  $\omega_s$  and  $s$  for the selection effect.

### **3 The instability model: parameter estimation and reconstruction of genome losses**

Our estimation strategy rests on the assumption that the large number of typed markers allows us to gather enough information on the sporadic loss process that its parameters  $\eta, \delta$ , and  $\tau$  can be estimated separately from  $\omega$  and the TSG location  $s$ . This is certainly the case when one is typing SNPs at high density genomewide,  $\eta, \delta$ , and  $\tau$  are constants across chromosomes and individuals, and there is only one TSG per chromosome—which is the situation we consider in this paper. However, the same assumption may be appropriate when  $\eta, \delta$ , and  $\tau$  vary across chromosomes and individuals, and when there are multiple TSGs—depending on the marker density and the value of the parameters. In the cases in which this assumption appears inadequate, one would not resort to the two stage strategy for estimating  $\eta, \delta, \tau, \omega$  and  $s$  that we describe below, but would need to simultaneously estimate all of these parameters. This is not difficult theoretically, and the likelihood derivatives given in the following can be used to describe such a maximization routine. We did not pursue this strategy because the nature of the data we were interested in made it unnecessary and computationally very intensive. Note that if a smaller number of markers is typed—so that the proposed assumption may be inadequate—the number of computations will also

be significantly reduced, making the simultaneous estimation strategy more feasible (simultaneous estimation is carried out, for example, in Newton *et al.* 1998).

To estimate the parameters of the instability model,  $\delta$ ,  $\eta$  and  $\tau$ , we use a maximum likelihood approach and a gradient algorithm. As anticipated, we consider a likelihood that is based solely on the instability component of our model:

$$\log \mathcal{L}_I(Y^1, \dots, Y^T | \eta, \delta, \tau) = \sum_{t=1}^T \sum_{k=1}^{22} \log P_I(X_k^t).$$

To describe the form of the first derivatives of the loglikelihood with respect to the three parameters of interest, it is easier to focus initially on one term  $P_I(X_k^t) = P_I(X)$ . When the hidden state  $\Pi = (\pi_i)$  is known, we obtain

$$P_I(X, \Pi) = P_I(x_1, \dots, x_m, \pi_1, \dots, \pi_m) = p(\pi_1) \prod_{i=1}^{m-1} t(\pi_{i+1} | \pi_i) \prod_{i=1}^m e(x_i | \pi_i).$$

Therefore,

$$\frac{\partial P_I(X, \Pi)}{\partial t(\pi_{i+1} | \pi_i)} = \frac{P_I(X, \Pi)}{t(\pi_{i+1} | \pi_i)} \quad \frac{\partial P_I(X, \Pi)}{\partial e(x_i | \pi_i)} = \frac{P_I(X, \Pi)}{e(x_i | \pi_i)}.$$

Now, recalling that the likelihood of the data is  $P_I(X) = \sum_{\pi_1, \dots, \pi_m} P_I(X, \Pi)$ , and that we can carry out the summations with respect to  $\pi_i$  using the forward and backward recursions, we obtain

the following expression for the partial derivative for  $\eta$ :

$$\begin{aligned}
\frac{\partial P_I(X)}{\partial \eta} &= \frac{\partial \sum_{\pi_1, \dots, \pi_m} P_I(X)}{\partial \eta} = \frac{\sum_{\pi_1, \dots, \pi_m} \partial P_I(X, \Pi)}{\partial \eta} \\
&= \sum_{\pi_1, \dots, \pi_m} \sum_{i=1}^{m-1} \frac{P_I(X, \Pi)}{t(\pi_{i+1} | \pi_i)} \frac{\partial t(\pi_{i+1} | \pi_i)}{\partial \eta} \\
&= \sum_{i=1}^{m-1} \sum_{\pi_i} \sum_{\pi_{i+1}} \frac{1}{t(\pi_{i+1} | \pi_i)} \frac{\partial t(\pi_{i+1} | \pi_i)}{\partial \eta} \sum_{\substack{\pi_1, \dots, \pi_{i-1} \\ \pi_{i+2}, \dots, \pi_m}} P_I(X, \Pi) \\
&= \sum_{i=1}^{m-1} \sum_{\pi_i} \sum_{\pi_{i+1}} \frac{P_I(X, \pi_i, \pi_{i+1})}{t(\pi_{i+1} | \pi_i)} \frac{\partial t(\pi_{i+1} | \pi_i)}{\partial \eta} \\
&= \sum_{i=1}^{m-1} \sum_{\pi_i} \sum_{\pi_{i+1}} \alpha(\pi_i) e(x_i | \pi_{i+1}) \beta(\pi_{i+1}) \frac{\partial t(\pi_{i+1} | \pi_i)}{\partial \eta},
\end{aligned}$$

where  $\alpha(\pi_i)$  and  $\beta(\pi_i)$  are computed by the forward and backward algorithms. Similarly, we have

$$\frac{\partial P_I(X)}{\partial \delta} = \sum_{\pi_1} \alpha(\pi_1) \beta(\pi_1) \frac{\partial t(\pi_1)}{\partial \delta} + \sum_{i=1}^{m-1} \sum_{\pi_i} \sum_{\pi_{i+1}} \alpha(\pi_i) e(x_i | \pi_{i+1}) \beta(\pi_{i+1}) \frac{\partial t(\pi_{i+1} | \pi_i)}{\partial \delta},$$

where  $t(\pi_1 = 1) = \delta$  and  $t(\pi_1 = 0) = 1 - \delta$ . And finally, for the parameter  $\tau$ :

$$\frac{\partial P_I(X)}{\partial \tau} = \sum_{i=1}^m \frac{\alpha(\pi_i = 1) \beta(\pi_i = 1)}{e(x_i | \pi_i = 1)} \frac{\partial e(x_i | \pi_i = 1)}{\partial \tau}.$$

The derivatives of the log-likelihood can be easily computed from the expressions given above and used to set up a gradient algorithm:

$$\begin{aligned}
\delta^{(t+1)} &= \delta^{(t)} + \lambda_\delta \frac{\partial \log P_I(X)}{\partial \delta} \Big|_{\delta=\delta^{(t)}} \\
\eta^{(t+1)} &= \eta^{(t)} + \lambda_\eta \frac{\partial \log P_I(X)}{\partial \eta} \Big|_{\eta=\eta^{(t)}} \\
\tau^{(t+1)} &= \tau^{(t)} + \lambda_\tau \frac{\partial \log P_I(X)}{\partial \tau} \Big|_{\tau=\tau^{(t)}},
\end{aligned}$$

with  $\lambda$  indicating the step size. It is also quite clear how different assumptions on variability of the parameters  $\eta, \delta, \tau$  across chromosomes and individuals will result in derivatives of the log-likelihood obtained using summations across different index sets.

Once an estimate of the instability parameters is obtained, this can be used to reconstruct the most likely genomic aberration profile  $\Pi_k^t$  for each of the individuals and chromosomes in the sample, using a standard Viterbi algorithm (see, Durbin *et al.*, 1998). This represents by itself an interesting output of our procedure, as it allows scientists to gather information on location and size of genomic losses from data that consist only in genotypes of cancer cells. Indeed, some researchers may consider this as the only output of interest, and not subscribe to the selection model for identification of TSGs that we will describe in the following section. For this reason, we preferred to opt for the 2-stage procedure.

## 4 The selection effect: likelihood ratio test to identify tumor suppressor gene locations

For the purpose of identifying the location  $s$  of a tumor suppressor gene and the increased probability of genomic loss  $\omega_s$  associated with this location, we conduct a series of likelihood ratio tests, where  $s$  is allowed to vary position across the entire chromosome, and the hypothesis  $H_0^s : \delta = \omega_s$  vs  $H_1^s : \delta \leq \omega_s$  is tested at each of the examined locations  $s$ . The genomic region for which the hypothesis  $H_0^s$  cannot be rejected will be considered as likely regions to harbor tumor suppressor genes, and the corresponding  $\omega_s$  will describe the selection effects. At this stage, the parameters  $\delta, \eta, \tau$  are considered known, so that for each explored location  $s$  we have to maximize the likelihood only with respect to the parameter  $\omega_s$ . Recall that for a given location  $s$ , the likelihood of the chromosome data is:

$$\mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega_s, \tau, s) = \prod_{t=1}^T (\omega_s P_I(X^t | \pi_s^t = 1) + (1 - \omega_s) P_I(X^t | \pi_s^t = 0)).$$

Notice that  $\mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega_s = \delta, \tau, s) = \prod_{t=1}^T P_I(X^t | \eta, \delta, \tau)$ . Furthermore, taking the logarithm of  $\mathcal{L}$  we obtain a concave function of  $\omega_s$ ; indeed  $\log \mathcal{L}(\omega_s) = \sum_{t=1}^T \log(\omega_s (P_I(X^t | \pi_s^t = 1) - P_I(X^t | \pi_s^t = 0)) + P_I(X^t | \pi_s^t = 0))$  is the sum of logarithms of affine functions, hence the sum

of concave functions is thus concave. This allows us to conclude that there is only one maximum for  $\log \mathcal{L}(\omega_s)$  and for  $\mathcal{L}(\omega_s)$ . In carrying out this series of likelihood ratio tests we follow the convention used in linkage genome scans of recording the logarithm base 10 of the inverse of the likelihood ratio, called LOD score; precisely, at each examined location  $s$  our test statistics  $L_s$  will be

$$L_s = \begin{cases} \log_{10} \frac{\mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega_s^*, \tau, s)}{\prod_{t=1}^T P_I(X^t | \eta, \delta, \tau)} & \omega_s^* = \operatorname{argmax} \mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega_s, \tau, s) > \delta \\ 0 & \omega_s^* < \delta \end{cases}$$

In terms perhaps more familiar to statisticians, that  $L_s$  can be interpreted as a profile log-likelihood ratio. The  $s^*$  location that maximizes  $L_s$  represents a candidate tumor suppressor gene location. To decide if the evidence in favor of  $\omega_s^* > \delta$  is sufficient, one needs to further examine the value of  $L_{s^*}$  and determine if such a difference in likelihood values is attributable to random chance or not. We will return to this point in the following.

To maximize  $\mathcal{L}(X^1, \dots, X^T | \delta, \eta, \omega_s, \tau, s)$  as a function of  $\omega_s$  it is convenient to use an EM algorithm. If we consider  $\pi_s^1, \dots, \pi_s^T$ , loss status at location  $s$  for the  $T$  cell-lines in the sample as missing data, we obtain the following complete data likelihood:

$$P(X^1, \dots, X^T, \pi_s^1, \dots, \pi_s^T | \omega, s) = \prod_{t=1}^T (\omega P_I(X^t | \pi_s^t = 1))^{\pi_s^t} ((1 - \omega) P_I(X^t | \pi_s^t = 0))^{(1 - \pi_s^t)},$$

leading to the complete data log-likelihood below:

$$\log P(X^1, \dots, X^T, \pi_s^1, \dots, \pi_s^T | \omega, s) = \sum_{t=1}^T \pi_s^t \log(\omega) + (T - \sum_{t=1}^T \pi_s^t) \log(1 - \omega) + \sum_{t=1}^T \log P_I(X^t | \pi_s^t).$$

The last term does not depend on  $\omega$ , so we can omit it from consideration in the following. The EM iterations will be based on the following expected complete data log-likelihood:

$$Q(\omega | \omega^{(l)}) = \sum_{t=1}^T E(\pi_s^t | X^t, \omega^{(l)}) \log(\omega) + (T - \sum_{t=1}^T E(\pi_s^t | X^t, \omega^{(l)})) \log(1 - \omega).$$

For the expectation step, we will have:

$$E(\pi_s^t | X^t, \omega^{(l)}) = \Pr(\pi_s^t = 1 | X^t, \omega^{(l)}) = \frac{P_I(X^t | \pi_s^t = 1) \omega^{(l)}}{P_I(X^t | \pi_s^t = 1) \omega^{(l)} + P_I(X^t | \pi_s^t = 0) (1 - \omega^{(l)})}.$$

Maximizing  $Q(\omega|\omega^{(l)})$  leads to  $\omega^{(l+1)} = \sum_{t=1}^T E(\pi_s^t|X^t, \omega^{(l)})/T$ .

We now return to the problem of determining which locations  $s$ , corresponding to high values of  $L_s$ , should be considered serious candidates for a TSG. First, note that  $L_s$  has a reasonable interpretation in terms of how much more likely the data is under  $H_1^s$  than  $H_0^s$ , and that the researcher may want to select a threshold value that better represents their interpretation of the study results. On purely statistical grounds, the determination of an appropriate cut off depends on the distribution of  $L_s$  under the null  $H_0^s$  and on the necessity of taking into account that multiple tests are being performed. Furthermore, notice that the tests  $L_{s_1}$  and  $L_{s_2}$  corresponding to two locations on the same chromosome are not independent. To determine a significance cut-off one ideally would like to know the distribution of the entire process  $\{L_s\}_s$  under the complete null hypothesis. Unfortunately, this is unknown at this stage. The marginal distribution of  $L_s$ , as  $T \rightarrow \infty$ , can be roughly approximated using the known results for likelihood ratio tests: under  $H_0^s$ ,  $2 \ln 10 L_s$  is asymptotically distributed as a 50:50 mixture of a mass at zero and  $\chi_{(1)}^2$  (the mass at zero derives from the fact that we place a constraint on the values of  $\omega > \delta$ , and the 0.5 mixing coefficient can be derived from the consistency and gaussianity of the MLE of  $\omega_s$ ). While this approximation of the distribution of  $L_s$  is rather crude, it provides us a guideline of what a reasonable significance cut-off may be. The appropriate cut-off for  $L_s$  depends on the distribution of  $L_s$  and, roughly speaking, on the number of “effectively independent” tests, which is determined by the length of the segment of the genome studied and the value of the  $\eta$  parameter. We suggest that once the instability parameters are estimated, a small scale simulation study be conducted where genotype data with the same structure as the real one is generated from the instability model, with no selection effect, and a cut-off for  $L_s$  that controls the desired measure of error rate to be determined. It may be of use to refer once again to the analogy with linkage mapping which carries through in terms of distribution for  $L_s$ : in these genetic mapping studies, a value of  $L_s$  greater than 3, or 3.5 is typically considered strong evidence in favor of  $H_1^s$  (Lander and Kruglyak, 1995).

## 5 Simulations and data example

In this section we illustrate the proposed methodology with a small scale simulation study and by applying it to a real dataset. When we simulate data, we consider the same structure of the Affymetrix 10k mapping arrays: 10913 SNPs covering 22 autosomal chromosomes with an average distance of 230 kb. To generate genotypes we use the population allele frequencies provided by Affymetrix. Our simulations are small scale, in that we did not consider a series of possible values for the model parameters, but we simply chose one and evaluated the performance of our model in one case. Our intention is not to investigate the performance of our model overall, but simply to illustrate its potential. We chose instability parameter values that are realistic for at least some biological settings:  $\delta = 0.1$ ,  $\eta = 0.2$ ,  $\kappa = 0.08$ ,  $\tau = 0.13$ . As far as the selection component of our model, we postulated one tumor suppressor gene at location 49.10 Mb on chromosome 1.

The first goal of our simulation was to evaluate the effectiveness of our procedure for the estimation of the instability parameters. We attempted estimation of these parameters using data from one individual only. To ensure that the presence of a TSG would not introduce distortions, we constrained the location 49.10 Mb on chromosome 1 to be lost and generated the rest of the data from the instability model. We repeated this 100 times and estimated in each case the three parameters  $\delta$ ,  $\eta$ , and  $\tau$ . Results are presented in Figure 1: mean and median of the estimators are both concordant with the true parameter values, and the spread is reasonably small.

Using the estimates of the instability parameters and the Viterbi algorithm, we subsequently reconstructed the loss status for each of the simulated cancer cell lines. We evaluated this reconstruction in terms of sensitivity (fraction of SNPs in regions of genomic loss that are correctly described) and specificity (fraction of SNPs in normal genomic regions that are correctly described). Figure 2 presents the histograms of sensitivity and specificity across the 100 cases: the performance is quite satisfactory, with average sensitivity 0.86 and specificity 0.99.

We then turn to investigate the effectiveness of our algorithm for the localization of the tumor

suppressor gene. For this problem we considered the same parameter settings described above,  $\omega=0.3$ , and genotypes from 50 cancer cell lines. We generated 100 datasets, and applied to each of them the two-step procedure consisting of estimating the instability parameters first and then attempting localization of the TSG. Our estimates of the instability parameters are quite good (unsurprisingly better than the one described above, due to the larger abundance of data). We evaluated the LOD curve at roughly 10,000 locations in the genome corresponding to the midpoint between each of the markers available. To determine the appropriate significance threshold for the LOD score curve, we conducted a simulation study using the true values of the instability parameters. Given the quality of our estimates, this does not represent a significant short-cut: the value 3.5 resulted in control of the family wise error rate at the 0.05 level. Figures 3 and 4 illustrate the results in one of our simulations. The LOD score curve across the entire genome is presented in Figure 3: only the area on chromosome one corresponding to the location of the TSG leads to values higher than 3.5. Figure 4 provides a more detailed illustration of the data on chromosome one and our analysis for this simulation: comparing panels (a) and (b) one can appreciate the denoising achieved with the reconstruction of the most likely loss profile. Simultaneously, the  $\omega_s$  estimates and the LOD curve (based on the simultaneous analysis of all cell lined) detect signals that are not enriched sufficiently at the level of individual tumors to be interpreted as losses.

When we analyze the results of the 100 simulations, we obtain a power to detect the TSG on chromosome 1 equal to 0.5. The actual FWER is 0.05 and the FDR 0.037 (to estimate these, we considered as a false discovery any location leading to a  $\text{LOD} > 3.5$  and more than  $1/\eta$  away from the true TSG). The relatively low power can be explained considering that  $\omega$  and the average size of losses (as determined by  $\eta$ ) are both relatively small. Furthermore, it is noteworthy that we generated our dataset setting  $\omega=0.3$ , but we did not fix the actual fraction of cancer cells in the sample to be 0.3; this reduces the overall power. Figure 5 illustrates how the LOD score at the TSG location varies as a function of the actual proportion of losses in the sample. From Figure 5 is also evident that larger  $\omega$  values will increase power; a larger sample size will also obviously

do so. The effect of  $\eta$  and  $\delta$  on power cannot be deduced from our simulations. Higher values of  $\delta$  would make sporadic losses more common, and higher values of  $\eta$  would make sporadic losses shorter (and hence less likely to overlap): changes in these parameters will affect the significance cut-off level and likely the power.

While the results of this small scale simulation study appear fairly positive, one must recall that we generated the data according to the model we used to analyze them. Our simulation allows us to verify the accuracy of our algorithm and evaluate the power of our procedure in an ideal situation, but it does not reflect the difficulties that we may encounter in the analysis of real data. To partially address this concern, we now turn to the analysis of a small dataset, collected by our colleagues at UCLA that motivated this methodology development. In a pilot study, Affymetrix 10k arrays were used to genotype 11 samples of primary brain tumors. In the initial analysis, whose results we report here, we assumed that the instability model parameters  $\delta$ ,  $\eta$ , and  $\tau$  were constant across chromosomes and cell lines. Their estimated values were  $\delta = 0.1638$ ,  $\eta = 0.6777$ ,  $\tau = 0.1398$ . A simulation study suggested the value of 3.8 as a cut-off for  $L_s$  in order to control FWER at 0.05. Only one region in the genome reached that level, on the short arm of chromosome 9: data and LOD score are reported in Figure 6. Interestingly, this region is known to harbor a tumor suppressor gene, referred to as p16 (Lucas *et al.*, 1995). Considering the data in Figure 6 leads to a few remarks. There are four samples that are almost entirely homozygous on the short arm of chromosome 9, and two with high loss rates: this is reflected in an overall higher level of the LOD curve. The maximum LOD score, however, is not reached in correspondence of the known TSG. The length of these losses is quite higher than what is predicted by the value of  $\eta$ : while this makes the signal quite strong in this region, it can also be interpreted as a model failure; we are currently exploring more careful analysis of the data, allowing instability parameters to vary across cell lines and chromosomes.

## 6 Discussion

We describe a model to reconstruct genomic losses in cancer cell lines from genotype data and to identify the likely location of tumor suppressor genes on the basis of genotype data from a collection of cancer cell lines. Our model for the process of genomic loss is taken from Newton *et al.* (1998), and we coupled it with a hidden Markov model to specify how the process of genomic loss is translated in observations on genotypes of the cancer cell line only. The study design we consider and the choices at the foundation of our computational strategy are based on the use of high density SNP genotyping, as obtained with Affymetrix 10k arrays, which allow scoring of 10,000 polymorphisms across the human genome.

When only genotype data from cancer cell lines are available, information on genomic losses is contained in stretches of homozygous markers: the longer the stretch of homozygosity, the less likely that such multilocus genotypes may be observed in absence of genomic aberrations. While in our model genomic loss is the only cause for increased multilocus homozygosity, there are other mechanisms that can potentially generate genotypes with similar characteristics: in particular, inbreeding (Leutenegger *et al.*, 2003) and linkage disequilibrium (Sabatti and Risch, 2002; Rosenberg and Calabrese, 2004) are important alternative explanations for increased homozygosity. We now briefly discuss the implications of these observations.

Linkage disequilibrium is the terminology used in genetics to indicate association between alleles at nearby markers. The model for genotypes given the genomic loss status  $\Pi$  that we proposed assumes linkage equilibrium, that is absence of association. As the distance between neighboring markers decreases, this assumption becomes inadequate and leads to identification of more losses than are really present in the data set. For data obtained with the Affymetrix genotyping array and described in this paper, this does not present a very serious problem: linkage disequilibrium, when present, typically extends at most across a window of 3-4 markers—which is generally too small to lead to the suspicion of a genomic loss. However, when the density of the scored polymorphisms

increases, it becomes necessary to incorporate the effect of linkage disequilibrium in our model. Indeed, it is possible to modify the emission probabilities of the HMM in order to account for the association between alleles at neighboring markers. This leads to a likelihood that cannot be described in the typical framework of HMM, but can still be evaluated with recursive formulas. We describe this model in a companion paper (Wang H. *et al.*, 2005).

The effect of inbreeding and genomic losses are more difficult to distinguish. Indeed, Leutenegger *et al.* (2003) have a model to estimate inbreeding coefficients that is practically identical to the HMM we use to describe genome instability. In this regard, the genomic loss profile that we reconstruct for each cell line under the instability model is valid only under the assumption of no inbreeding (which, incidentally, is likely to hold in the type of population samples collected for LOH studies). The instability component of our model, instead, captures an effect that is clearly distinct from inbreeding: the propensity of multiple cell lines to share losses in the same region. The inference on TSG location, hence, should be robust to the presence of inbreeding when sufficient numbers of independent cancers are analyzed.

Finally, we note that there are some extensions to our model that can be easily carried through and that we did not implement in the present paper partly because of time constraints and partly because they did not appear necessary to perform the analysis of the data we set out to study, namely genotypes obtained with the Affymetrix 10k array. One such extension involves the introduction of genotyping error parameter in the emission probabilities. This parameter can either be assumed known or estimated from the data (however, estimating it from another dataset is likely to produce better results). Similarly, the  $\kappa$  parameter can be varied across markers or cell lines. Another extension consists in the simultaneous estimation of all the parameters  $\eta, \delta, \tau, \omega_s$ , which, as we have discussed in the previous section, may be important for datasets that do not provide genome-wide genotyping. Furthermore, while we described our model assuming the instability parameters held constant across cell-lines and chromosomes, one can easily relax this assumption.

## Acknowledgments

Chiara Sabatti and Hui Wang were partially supported by NSF grant DMS0239427 and NIH/NIDOCDC grant DC04224. Chiara Sabatti also acknowledge support from ASA/Ames grant NCC2-1364 and USPHS grant GM53275. DNA microarray data were generated with support from NCI grant U01 CA88127.

## References

- Bignell, G., Huang, J., Greshock, J., Watt, S., Butler, A., West, S., Grigorova, M., Jones, K.W., Wei, W., Stratton, M., Futreal, P., Weber, B., Shapero, M., Wooster, R. (2004) "High-resolution analysis of DNA copy number using oligonucleotide microarrays," *Genome Res.* 14: 287–95.
- Cox, C., Bignell, G., Greenman, C., Stabenau, A., Warren, W., Stephens, P., Davies, H., Watt, S., Teague, J., Edkins, S., Birney, E., Easton, D., Wooster, R., Futreal, P., Stratton, M. (2005) "A survey of homozygous deletions in human cancer genomes," *PNAS* 102: 4542–4547.
- Durbin, E., Eddy, S., Krogh, A., and Mitchinson, G. (1999) *Biological sequence analysis*, Cambridge University Press.
- Fridlyand, J., Snijders, A., Pinkel, D., Albertson, D., and Jain, A. (2004) "Hidden Markov models approach to the analysis of array CGH data," *Journal of Multivariate Analysis* 90: 132–153.
- Huang, J., Wei, W., Zhang, J., Liu, G., Bignell, G., Stratton, M., Futreal, P., Wooster, R., Jones, K., Shapero, M. (2004) "Whole genome DNA copy number changes identified by high density oligonucleotide arrays," *Hum. Genomics* 1: 287–99.
- Iafate, A., Feuk, L., Rivera, M., Listewnik, M., Donahoe, P., Qi, Y., Scherer, SW, Lee, C. (2004) "Detection of large-scale variation in the human genome," *Nature Genetics* 36: 949–51.

- Lai, Y., Zhao, H. (2005) "A statistical method to detect chromosomal regions with DNA copy number alterations using SNP-array-based CGH data," *Comput. Biol. Chem.* 29: 47–54.
- Lander, E., and Kruglyak, L. (1995) "Genetic dissection of complex traits: guidelines for interpreting and reporting linkage results," *Nature Genetics* 11: 241–247.
- Leutenegger, A., Prum, B., Genin, E., Verny, C., Lemainque, A., Clerget-Darpoux, F., Thompson, E. (2003) "Estimation of the inbreeding coefficient through use of genomic data," *Am. J. Hum. Genet.* 73:516–23.
- Lin, M., Wei, L., Sellers, W., Lieberfarb, M., Wong, W., Li, C.(2004) "dChipSNP: significance curve and clustering of SNP-array-based loss-of-heterozygosity data," *Bioinformatics* 20: 1233–40.
- Lukas, J., Parry, D., Aagaard, L., Mann, D. J., Bartkova, J., Strauss, M., Peters, G., Bartek, J. (1995) "Retinoblastoma-protein-dependent cell-cycle inhibition by the tumour suppressor p16," *Nature* 375: 503–506.
- Miller, B., Wang, D., Krahe, R., Wright, F. (2003) "Pooled analysis of loss of heterozygosity in breast cancer: a genome scan provides comparative evidence for multiple tumor suppressors and identifies novel candidate regions," *Am. J. Hum. Genet.* 73: 748–67.
- Newton, M., Lee, Y. (2000) "Inferring the location and effect of tumor suppressor genes by instability-selection modeling of allelic-loss data," *Biometrics* 56: 1088–97.
- Newton, M., Gould M., Reznikoff, C., Haag, J. (1998) "On the statistical analysis of allelic-loss data," *Stat. Med.* 17: 1425–45.
- M.A. Newton (2002) "Discovering combinations of genomic alterations associated with cancer," *Journal of the American Statistical Association* 97: 931–942.

- Olshen, A., Venkatraman, E., Lucito, R., Wigler, M. (2004) “Circular binary segmentation for the analysis of array-based DNA copy number data,” *Biostatistics* 5: 557–72.
- Pinkel, D., Se Graves, R., Sudar, D., Clark, S., Poole, I., Kowbel, D., Collins, C., Kuo, W., Chen, C., Zhai, Y., Dairkee, S., Ljung, B., Gray, J., Albertson, D. (1998) “High resolution analysis of DNA copy number variation using comparative genomic hybridization to microarrays,” *Nature Genetics* 20: 207–11.
- Pollack, J., Sorlie, T., Perou, C., Rees, C., Jeffrey, S., Lonning, P., Tibshirani, R., Botstein, D., Borresen-Dale, A., and Brown, P. (2002) “Microarray analysis reveals a major direct role of DNA copy number alteration in the transcriptional program of human breast tumors,” *PNAS* 99: 12963–12968.
- Rosenberg, N., Calabrese, P. (2004) “Polyploid and multilocus extensions of the Wahlund inequality,” *Theoretical Population Biology* 66: 381–391.
- Sabatti, C. and Risch, N. (2002) “Homozygosity and linkage disequilibrium,” *Genetics* 160: 1707–1719.
- Sebat, J., Lakshmi, B., Troge, J., Alexander, J., Young, J., Lundin, P., Maner, S., Massa, H., Walker, M., Chi, M., Navin, N., Lucito, R., Healy, J., Hicks, J., Ye, K., Reiner, A., Gilliam, T., Trask, B., Patterson, N., Zetterberg, A., Wigler, M. (2004) “Large-scale copy number polymorphism in the human genome,” *Science* 305: 525–8.
- Wang, H., Service, S., Freimer, N., and Sabatti, C. (2005) “Detecting large scale genomic variation through high density SNP genotyping,” *manuscript in preparation*.
- Wang, P., Kim, Y., Pollack, J., Narasimhan, B., Tibshirani, R. (2005) “A method for calling gains and losses in array CGH data,” *Biostatistics* 6: 45–58.

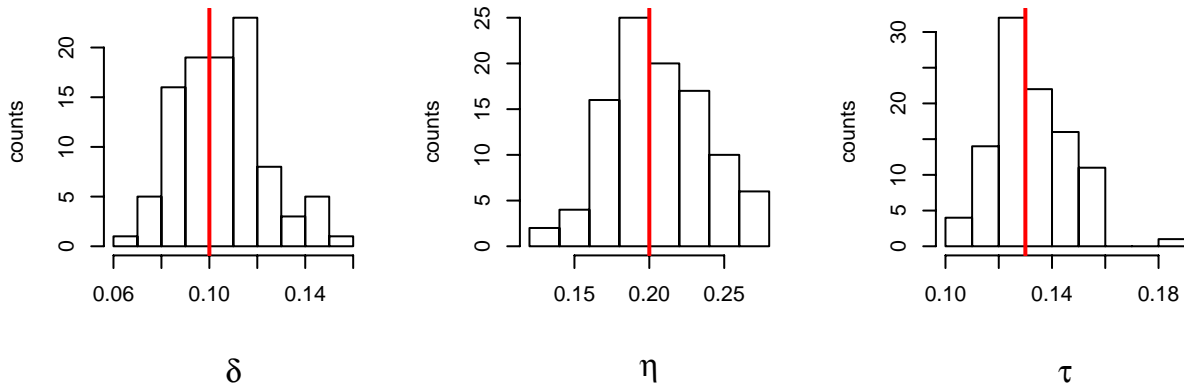


Figure 1: **Instability Parameters Reconstruction.** Histograms of the estimated values for each of the instability parameters using genotypes from one cell line in 100 simulations. The vertical red lines indicate the true value of the parameters.

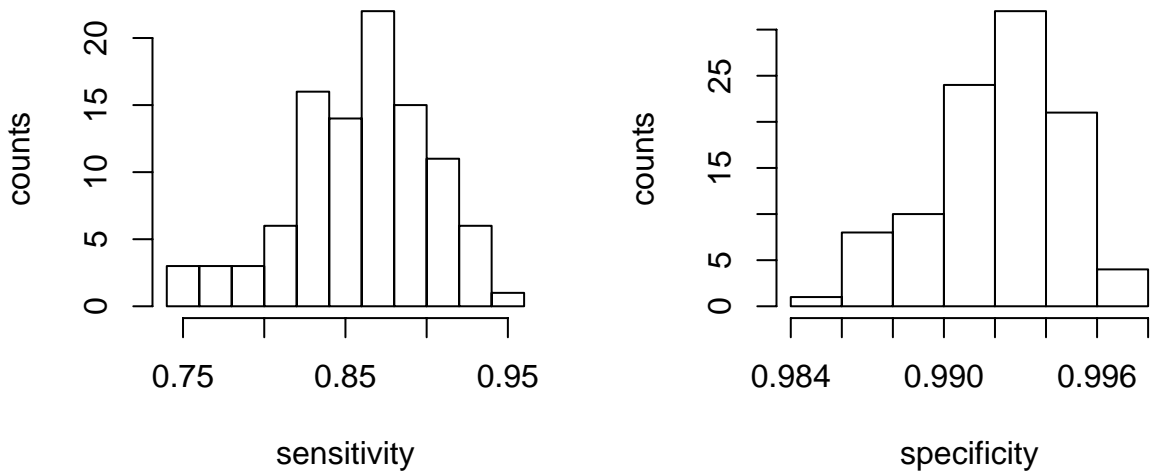


Figure 2: **Reconstruction of Loss Status.** Histograms of sensitivity and specificity of the loss process reconstruction using instability model in 100 simulations.

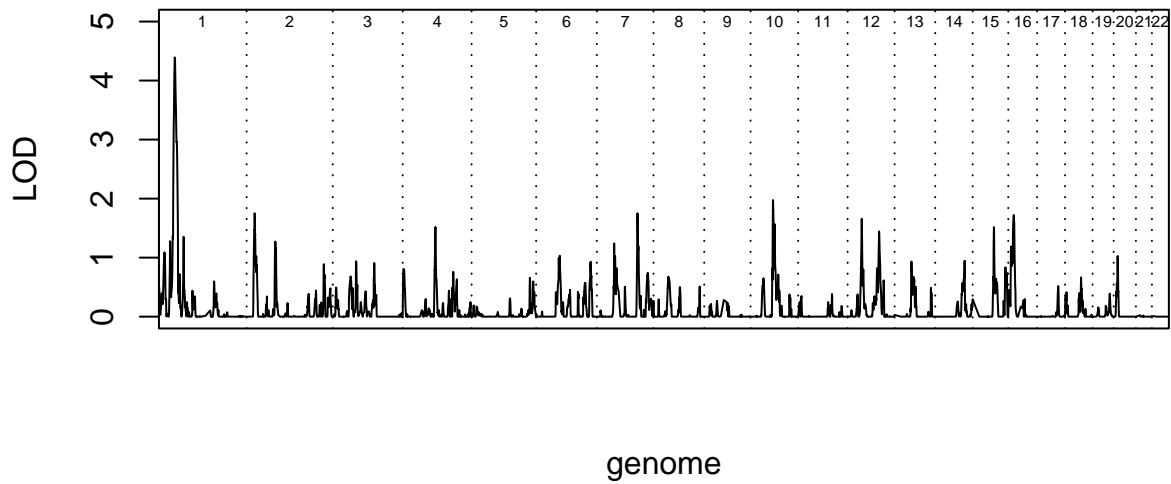


Figure 3: **LOD score for TSG.** Profile log-likelihood ratio for the location  $s$  of the tumor suppressor gene across the genome for one case of the simulation.

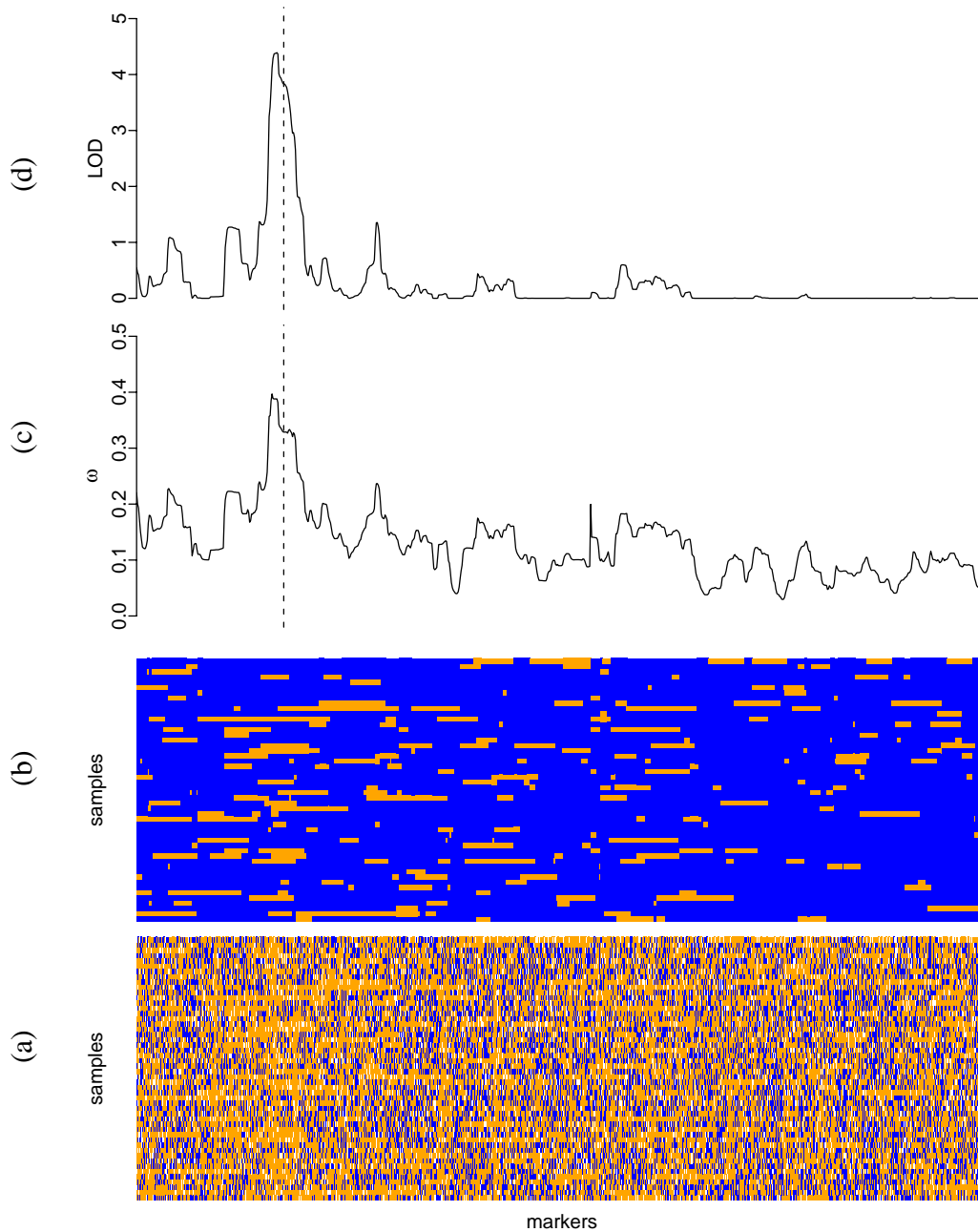


Figure 4: **Localization of the TSG on Chromosome One.** Panel (a) represents the data generated by our simulation and used in the reconstruction. Each row represents an individual and each column a SNP. Polymorphisms are ordered according to their genomic position. Heterozygous genotypes are colored in blue, homozygous genotypes in orange, and missing calls are left blank. For ease of display, inter-marker distances are depicted as constant. Panel (b) displays the corresponding loss profiles reconstructed with the instability model. Panel (c) reports the estimated value for  $\omega_s$  for each of the genomic positions. Finally, in panel (d), we present the LOD score curve for  $s$ . The true location of the TSG is indicated with a dashed line.

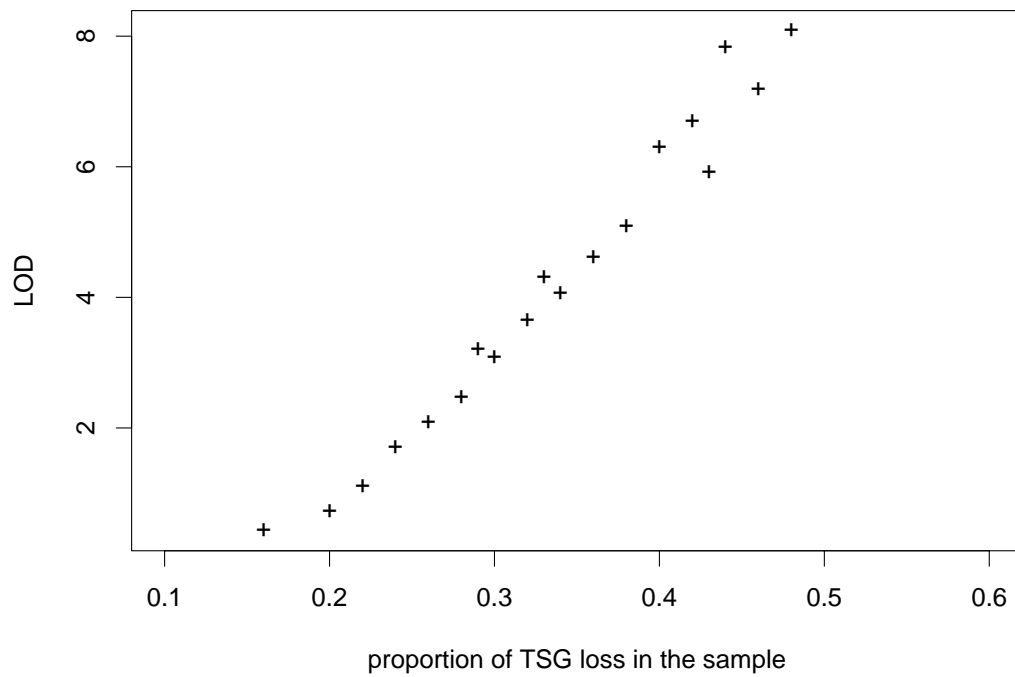


Figure 5: **LOD value at the TSG.** Simulated datasets were grouped according to the actual proportion of samples with a loss in the TSG locations. The average LOD score at the TSG was computed within these groups and is plotted against the corresponding proportion of losses.

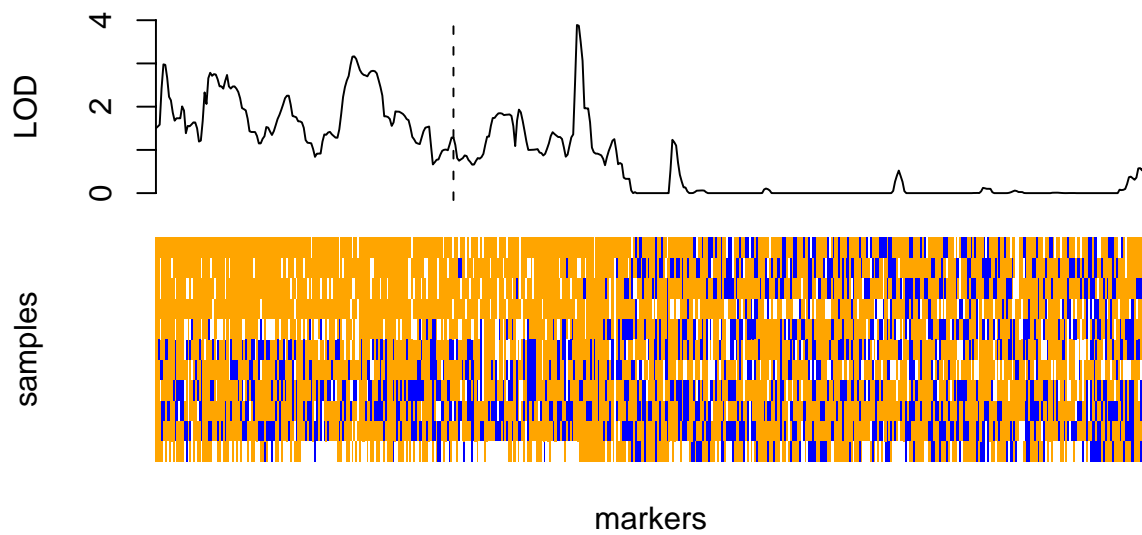


Figure 6: **Data Analysis.** We present the data (coded in the same format as in Figure 4) and the LOD curve for the localization of a tumor suppressor gene relative to chromosome 9. The location of the known TSG p16 is indicated with a dashed line.