

# Note

## Fast and Accurate Estimation of the Population-Scaled Mutation Rate, $\theta$ , From Microsatellite Genotype Data

Arindam RoyChoudhury<sup>1</sup> and Matthew Stephens

Department of Statistics, University of Washington, Seattle, Washington 98195-4322

Manuscript received August 2, 2005

Accepted for publication February 24, 2007

### ABSTRACT

We present a new approach for estimation of the population-scaled mutation rate,  $\theta$ , from microsatellite genotype data, using the recently introduced “product of approximate conditionals” framework. Comparisons with other methods on simulated data demonstrate that this new approach is attractive in terms of both accuracy and speed of computation. Our simulation experiments also demonstrate that, despite the theoretical advantages of full-likelihood-based methods, methods based on certain summary statistics (specifically, the sample homozygosity) can perform very competitively in practice.

PATTERNS of genetic variation in population samples contain important information on both the biological mechanisms (*e.g.*, mutation, recombination, gene conversion, selection) and aspects of population demographic history (*e.g.*, population expansions, bottlenecks, and migration rates). However, extracting this information is often tricky. The simplest methods are based on matching summaries of the data (*e.g.*, expected heterozygosity or average pairwise distances between alleles) to their expected values. Although these methods are attractive in their simplicity, summarizing the genotype data with a single number in this way risks losing information. More complex methods that use sophisticated computations to approximate the full likelihood of the data (GRIFFITHS and TAVARÉ 1994a,b; KUHNER *et al.* 1995; IORIO *et al.* 2005) are more efficient in principle, but typically are difficult to implement, and may take impractical amounts of time to produce reliable results (STEPHENS and DONNELLY 2000; FEARNHEAD and DONNELLY 2001). This has limited their usefulness in practice. Indeed, in some settings the computational complexities of full-likelihood-based approaches are so daunting that many researchers have turned to approximate methods (*e.g.*, HUDSON 2001; McVEAN *et al.* 2002; FEARNHEAD and DONNELLY 2002; LI and STEPHENS 2003), often with considerable success (*e.g.*, CRAWFORD *et al.* 2004; McVEAN *et al.* 2004). Thus far, applications of these approximate methods have been to data on single-

nucleotide polymorphisms (SNPs). Here we extend one of these methods, the PAC likelihood approach of LI and STEPHENS (2003), to estimate the scaled mutation parameter  $\theta$  ( $= 2N\mu$ , where  $N$  is the effective haploid population size and  $\mu$  is the mutation probability per meiosis) from microsatellite data. Simulation results suggest that this method is as accurate as full-likelihood-based approaches and considerably faster.

**Models and methods:** We consider a simple scenario, where we genotype a single microsatellite locus in  $n$  haploid individuals, or  $n/2$  diploid individuals, sampled from a random-mating population that has been evolving neutrally with constant (haploid) size  $N$  according to a Wright–Fisher model. Let  $a_1, \dots, a_n$  denote the observed alleles (number of repeats of the microsatellite motif). We assume that the locus evolves according to a symmetric stepwise mutation mechanism, where if a mutation occurs in a transmission then the offspring’s allele length increases or decreases (with equal probability) by one from the progenitor allele. Although this model is simplistic, it is widely used and is the basis for all the methods of estimating  $\theta$  that we consider here. However, our approach could be easily modified to deal with other mutation models (*e.g.*, those described in CALABRESE and DURRETT 2003).

There exist two broad categories of approach for estimating  $\theta$  in this context. The first is moment estimators based on summary statistics. KIMMEL *et al.* (1998) include two such estimators (their Equations 14 and 15). The first one, the homozygosity estimator, is given by

$$\hat{\theta}_H = 0.5 \left( \hat{P}_0^{-2} - 1 \right), \quad (1)$$

<sup>1</sup>Corresponding author: Wakeley Lab, 4092-4100 Biological Laboratories, 16 Divinity Ave., Harvard University, Cambridge, MA 02138.  
E-mail: aroy@fas.harvard.edu

where  $\hat{P}_0$  is an unbiased estimate of the population homozygosity,

$$\hat{P}_0 = \frac{n \sum_{k=1}^r p_k^2 - 1}{n - 1}, \tag{2}$$

where  $r$  is the number of different alleles found in the population, and  $p_i$  is the sample frequency of the  $i$ th allele. The second estimator is

$$\hat{\theta}_V = (2/(n - 1)) \sum_{i=1}^n (a_i - \bar{a})^2, \tag{3}$$

where  $\bar{a}$  is the mean of the  $a_i$ 's. The estimator  $\hat{\theta}_H$  is based on the limiting expected homozygosity in a continuous-time Wright–Fisher model, whereas  $\hat{\theta}_V$  is based on the limiting expected value of the within-population component of genetic variance in the same model KIMMEL *et al.* (1998).

The second category is full-likelihood-based approaches, including maximum-likelihood and Bayesian approaches, which base inference on the likelihood

$$L(\theta) = \Pr(a_1, \dots, a_n | \theta). \tag{4}$$

In principle full-likelihood-based approaches are more efficient than moment estimators based on summary statistics. However, they are considerably harder to implement because the likelihood (4) cannot be computed directly. Instead, the likelihood can be approximated using computational methods such as Markov chain Monte Carlo (MCMC) or importance sampling. WILSON and BALDING (1998) and BEERLI and FELSENSTEIN (2001) describe two such approaches. WILSON and BALDING (1998) take a Bayesian approach, specifying prior distributions for  $N$  and  $\mu$ , and use an MCMC scheme to draw samples from the posterior distribution of  $\theta$ . This method is implemented in the software MICSAT, which we downloaded from <http://www.maths.abdn.ac.uk/~ijw/downloads/download.htm>. BEERLI and FELSENSTEIN (2001) also use a (different) MCMC scheme; but instead of performing a Bayesian analysis, they use it to compute a likelihood surface for  $\theta$  (and also, in the case of samples from multiple populations, a set of migration rates among populations; however, here we deal with a sample from a single random-mating population, and so their approach can be used to estimate  $\theta$  alone). This method is implemented by the program Migrate (version 1.7.3), which we downloaded from <http://evolution.genetics.washington.edu/lamarc/migrate.html>.

In this article we take a different approach, following LI and STEPHENS (2003) who suggest approximating the likelihood (4) by exploiting the identity

$$\begin{aligned} &\Pr(a_1, \dots, a_n | \theta) \\ &= \Pr(a_1 | \theta) \Pr(a_2 | a_1; \theta) \dots \Pr(a_n | a_1, \dots, a_{n-1}; \theta). \end{aligned} \tag{5}$$

Although the conditional distributions on the right-hand side of this equation are unknown for most models of interest, they are amenable to approximation (*e.g.*, STEPHENS and DONNELLY 2000; FEARNHEAD and DONNELLY 2001; LI and STEPHENS 2003). Substituting such an approximation,  $\hat{\pi}$  say, into the right-hand side yields an approximate likelihood, which LI and STEPHENS (2003) term the “product of approximate conditionals” (PAC) likelihood:

$$L_{\text{PAC}}(\theta) = \hat{\pi}(a_1 | \theta) \hat{\pi}(a_2 | a_1; \theta) \dots \hat{\pi}(a_n | a_1, \dots, a_{n-1}; \theta). \tag{6}$$

LI and STEPHENS (2003) applied this idea to estimate recombination rates (but not mutation rates) from SNP data and showed the resulting estimates to be competitive with the best available methods for that problem.

Here we show that an analogous approach also works for estimating  $\theta$  from microsatellite data. For the conditional distributions  $\hat{\pi}(a_{k+1} | a_1, \dots, a_k; \theta)$  on the right-hand side of (6) we use the approximation suggested by STEPHENS and DONNELLY (2000). This approximation is based on the idea that the next sampled allele,  $a_b$ , will differ by a random number of mutations (which will typically be a small number of mutations and quite possibly 0 mutations) from a randomly chosen existing allele  $(a_1, \dots, a_k)$ . STEPHENS and DONNELLY (2000, p. 616) assume that the number of mutations,  $m$ , has a geometric distribution, with  $\Pr(m = 0) = k/(k + \theta)$ . The assumption of a geometric distribution is motivated by the fact that the resulting approximation is exact for the case  $k = 1$ ; and the assumption on  $\Pr(m = 0)$  is motivated by the fact that the resulting approximation is exact [and results in the well-known Ewens sampling formula (EWENS 1972)] for so-called “parent-independent mutation” (PIM) models, where the type of a mutant offspring is independent of the type of the progenitor allele. Of course, the stepwise mutation is not PIM, so the approximation is not exact in our setting. Part of our aim here is to show that the approximation is good enough to provide accurate estimates for  $\theta$ .

Mathematically, the approximation suggested by STEPHENS and DONNELLY (2000) is

$$\hat{\pi}(a_{k+1} | a_1, \dots, a_k; \theta) = (1/k) \sum_{i=1}^k \sum_{m=0}^{\infty} (1 - q_k) q_k^m (P^m)_{a_i a_{k+1}}, \tag{7}$$

where  $q_k = \theta/(k + \theta)$  and  $P$  is a mutation matrix, whose  $(i, j)$ th element is the probability that the type of an offspring is of type  $j$ , given that the progenitor is of type  $i$  and a mutation occurs. To ease comparison with other approaches, we assume a symmetric stepwise mutation mechanism, so that

$$P_{ij} = \begin{cases} 0.5, & \text{if } |i - j| = 1 \\ 0, & \text{otherwise.} \end{cases}$$

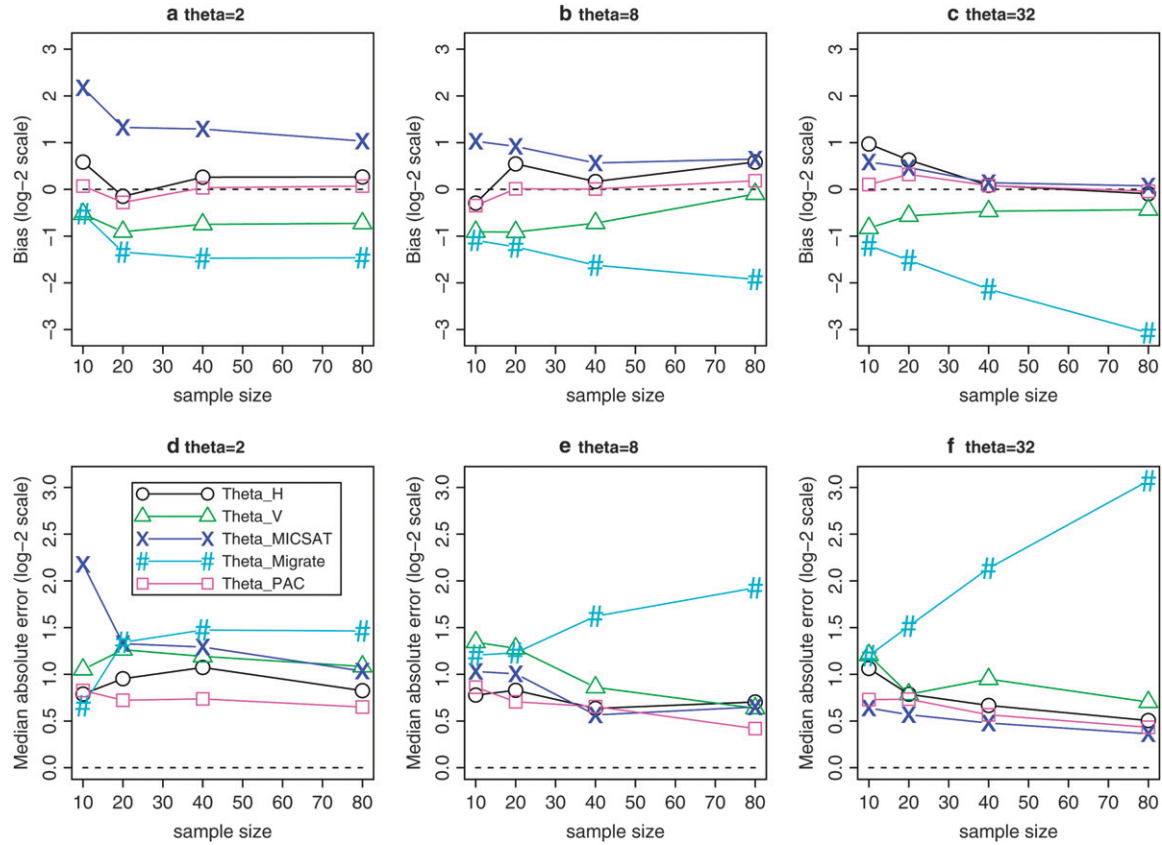


FIGURE 1.—Comparison of the “bias” (a–c) and “accuracy” (d–f) of different estimators. Each section has five curves, one for each estimator:  $\circ$ ,  $\hat{\theta}_H$ ;  $\triangle$ ,  $\hat{\theta}_V$ ;  $\times$ ,  $\hat{\theta}_{\text{MICSAT}}$ ;  $\#$ ,  $\hat{\theta}_{\text{Migrate}}$ ; and  $\square$ ,  $\hat{\theta}_{\text{PAC}}$ . In a–c the curves show the median value of  $\log \hat{\theta} - \log \theta$  for different haploid sample sizes  $n = 10, 20, 40$ , and  $80$ . In d–f each curve shows the median of  $|\log \hat{\theta} - \log \theta|$  for the same values of  $n$ . We used a coalescent-based simulation program, kindly provided by P. Fearnhead, to simulate samples of microsatellite alleles randomly sampled from a population evolving according to the Wright–Fisher model, with stepwise mutation. (This model underlies all the methods we compare here.) For each different  $\theta$ , and for each different  $n$ , we simulated 50 data sets. For each data set we estimated  $\theta$  using each of the methods and compared the estimated value of  $\theta$  with the true value of  $\theta$  used to generate the data. Approximate run times for a single data set of size  $n = 80$ , on a desktop computer with 3GHz CPU, were  $\sim 10$  min for MICSAT,  $\sim 45$  min for Migrate,  $\sim 10$  sec for our method, and  $< 1$  sec for the summary statistic methods.

We note that, unlike in STEPHENS and DONNELLY (2000), we do not impose any reflecting boundaries on the mutation process, although this would be straightforward to do. (Thus, the matrix  $P$  has infinitely many rows and columns.) It would also be straightforward to incorporate nonstepwise moves (*e.g.*, NIELSEN 1997) or indeed any other desired form for  $P$ .

This choice of  $P$  has the convenient, although not essential, property that the approximation (7) simplifies, to

$$\hat{\pi}(a_{k+1} | a_1, a_2, \dots, a_k; \theta) = (1/k) \sum_{i=1}^k \sqrt{\frac{1-q_k}{1+q_k}} \left( \frac{(1-\sqrt{1-q_k^2})/q_k}{q_k} \right)^{|a_i-a_{k+1}|} \tag{8}$$

This follows from rewriting (7) as

$$\hat{\pi}(a_{k+1} | a_1, \dots, a_k; \theta) = (1/k) \sum_{i=1}^k (1-q_k) (I - q_k P)_{a_i, a_{k+1}}^{-1} \tag{9}$$

and noting that the matrix with  $(i, j)$ th element

$$(I - qP)_{ij}^{-1} = \left( 1/\sqrt{1-q^2} \right) \left( \left( (1-\sqrt{1-q^2})/q \right)^{|i-j|} \right) \tag{10}$$

is the inverse of  $(I - qP)$ . Equation 10 can be verified by straightforward algebra, multiplying a row of  $(I - qP)$  by a column of  $(I - qP)^{-1}$  defined by (10).

Substituting (8) into (6) for  $k = 1, \dots, n - 1$  gives a PAC likelihood for this problem. Note that, as in LI and STEPHENS (2003), the resulting PAC likelihood is not invariant to the ordering of the sampled alleles  $a_1, a_2, \dots, a_n$ . To deal with this, we take the same approach as LI and STEPHENS (2003); we average (4) over 10 random permutations of  $a_1, a_2, \dots, a_n$ . [Results (not shown) obtained using a single random permutation were similar in accuracy.] We use  $\hat{\theta}_{\text{PAC}}$  to denote the value of  $\theta$  that maximizes this function [found numerically by computing  $L_{\text{PAC}}(\theta)$  on a dense grid of values for  $\theta$ ].

**Comparisons:** We compared the properties of our PAC-based estimator  $\hat{\theta}_{\text{PAC}}$  with other available methods described above: the moment-based estimators  $\hat{\theta}_{\text{H}}$  and  $\hat{\theta}_{\text{V}}$  and the full-likelihood-based estimators  $\hat{\theta}_{\text{MICSAT}}$  and  $\hat{\theta}_{\text{Migrate}}$ . To be precise,  $\hat{\theta}_{\text{MICSAT}}$  is the mean of 10,000 draws from the posterior distribution for  $\theta$  obtained using the program MICSAT with default parameter values, and  $\hat{\theta}_{\text{Migrate}}$  is the value of  $\theta$  that maximizes the approximate likelihood computed using Migrate, again with default parameter values.

Figure 1 compares “bias” (or, more accurately, median error) and “accuracy” (median absolute error) of the resulting estimates, on a log scale. Making comparisons on the log scale means that, for example, underestimating  $\theta$  by a factor of 2 is considered equally good—or bad—as overestimating by a factor of 2. We use medians rather than means because the means are infinite, due to the fact that there is a small finite probability of each estimator being 0 (and therefore giving a log of  $-\infty$ ); see also LI and STEPHENS (2003).

For the scenarios we consider,  $\hat{\theta}_{\text{PAC}}$ ,  $\hat{\theta}_{\text{H}}$ ,  $\hat{\theta}_{\text{HU}}$ , and  $\hat{\theta}_{\text{MICSAT}}$  are consistently better (smaller bias and smaller mean absolute error) than  $\hat{\theta}_{\text{V}}$  and  $\hat{\theta}_{\text{Migrate}}$ . If anything the results for  $\hat{\theta}_{\text{PAC}}$  seem very slightly better than the other three, especially for small values of  $\theta$  (according to a paired Wilcoxon signed-rank test, the improvement in accuracy over  $\hat{\theta}_{\text{MICSAT}}$  is significant at  $P < 0.05$  for all values of  $n$  considered at  $\theta = 2$  and for  $n = 10, 20, 80$  at  $\theta = 8$ ; the improvement over  $\hat{\theta}_{\text{H}}$  is significant at  $P < 0.05$  for all values of  $n$  considered at  $\theta = 2$ , for  $n = 20, 80$  at  $\theta = 8$ , and for  $n = 10, 20, 40$  at  $\theta = 32$ ). However, the differences may be too small to be practically important, and in some sense a direct comparison with  $\hat{\theta}_{\text{MICSAT}}$  is inappropriate, since it is based on a particular prior distribution for  $\theta$ .

One additional notable finding from our simulations is that, between the summary statistic estimators,  $\hat{\theta}_{\text{H}}$  performs considerably better than  $\hat{\theta}_{\text{V}}$ . Indeed, the finding that  $\hat{\theta}_{\text{H}}$  performs competitively with the likelihood-based methods is, as far as we are aware, novel. While we have no intuitive explanation for this good performance, the poor performance of  $\hat{\theta}_{\text{V}}$  might perhaps have been expected, for the following reason. Equation 3 for  $\hat{\theta}_{\text{V}}$  can be rewritten as  $\hat{\theta}_{\text{V}} = (1/n(n-1)) \sum_{i=1}^n \sum_{j=1}^n (a_i - a_j)^2$ . Thus  $\hat{\theta}_{\text{V}}$  is the mean squared pairwise difference between sampled microsatellite repeats. In the context of sequence data, the corresponding estimate for  $\theta$  (per base pair) is the mean pairwise distance (per base pair) between sampled haplotypes, also known as the nucleotide diversity, and this is known to be an inconsistent estimator for  $\theta$  in that context (*e.g.*, DONNELLY and TAVARÉ 1995).

We interpret the poorer performance of  $\hat{\theta}_{\text{Migrate}}$  as indicating that, even in this relatively simple setting, with only a single parameter to be estimated and no migration, the default run lengths we used were insufficient to

provide an accurate approximation to the maximum-likelihood estimates. In more complex settings, involving migration, for example, obtaining an accurate estimate of the likelihood surface, and the location of its maximum, seems likely to be still more challenging. Although some work would be necessary to extend our PAC-likelihood method to these settings, our results here, and in LI and STEPHENS (2003), suggest that this effort may be worthwhile.

We thank two anonymous referees for helpful comments on the submitted version of this manuscript. This work was supported by National Institutes of Health grant HG/LM02585 to M.S.

## LITERATURE CITED

- BEERLI, P., and J. FELSENSTEIN, 2001 Maximum likelihood estimation of a migration matrix and effective population sizes in  $n$  subpopulations by using a coalescent approach. *Proc. Natl. Acad. Sci. USA* **98**(8): 4563–4568.
- CALABRESE, P., and R. DURRETT, 2003 Dinucleotide repeats in the *Drosophila* and human genomes have complex, length-dependent mutation processes. *Mol. Biol. Evol.* **20**: 715–725.
- CRAWFORD, D., T. BHANGALE, N. LI, G. HELLENTHAL, M. RIEDER *et al.*, 2004 Evidence for substantial fine-scale variation in recombination rates across the human genome. *Nat. Genet.* **36**: 700–706.
- DONNELLY, P., and S. TAVARÉ, 1995 Coalescents and genealogical structure under neutrality. *Annu. Rev. Genet.* **29**: 401–421.
- EWENS, W. J., 1972 The sampling theory of selectively neutral alleles. *Theor. Popul. Biol.* **3**: 87–112.
- FEARNHEAD, P. N., and P. DONNELLY, 2001 Estimating recombination rates from population genetic data. *Genetics* **159**: 1299–1318.
- FEARNHEAD, P. N., and P. DONNELLY, 2002 Approximate likelihood methods for estimating local recombination rates. *J. R. Stat. Soc. Ser. B* **64**: 657–680.
- GRIFFITHS, R. C., and S. TAVARÉ, 1994a Ancestral inference in population genetics. *Stat. Sci.* **9**: 307–319.
- GRIFFITHS, R. C., and S. TAVARÉ, 1994b Simulating probability distributions in the coalescent. *Theor. Popul. Biol.* **46**: 131–159.
- HUDSON, R. R., 2001 Two-locus sampling distribution and their application. *Genetics* **159**: 1805–1817.
- IORIO, M. D., R. C. GRIFFITHS, R. LEBLOIS and F. ROUSSET, 2005 Stepwise mutation likelihood computation by sequential importance sampling in subdivided population models. *Theor. Popul. Biol.* **68**: 41–53.
- KIMMEL, M., R. CHAKRABORTY, J. P. KING, M. BAMSHAD, W. S. WATKINS *et al.*, 1998 Signatures of population expansion in microsatellite repeat data. *Genetics* **148**: 1921–1930.
- KUHNER, M. K., J. YAMATO and J. FELSENSTEIN, 1995 Estimating effective population size and mutation rate from sequence data using Metropolis–Hastings sampling. *Genetics* **140**: 1421–1430.
- LI, N., and M. STEPHENS, 2003 Modeling linkage disequilibrium, and identifying recombination hotspots using SNP data. *Genetics* **165**: 2213–2233.
- MCVEAN, G., P. AWADALLA and P. FEARNHEAD, 2002 A coalescent-based method for detecting and estimating recombination from gene sequences. *Genetics* **160**: 1231–1241.
- MCVEAN, G. A. T., S. R. MYERS, S. HUNT, P. DELOUKAS, D. R. BENTLEY *et al.*, 2004 The fine-scale structure of recombination rate variation in the human genome. *Science* **304**: 581–584.
- NIELSEN, R., 1997 A likelihood approach to population samples of microsatellite alleles. *Genetics* **146**: 711–716.
- STEPHENS, M., and P. DONNELLY, 2000 Inference in molecular population genetics. *J. R. Stat. Soc. Ser. B* **62**: 605–655.
- WILSON, I. J., and D. J. BALDING, 1998 Genealogical inference from microsatellite data. *Genetics* **150**: 499–510.