

Supporting Information

Genome-wide fine-scale recombination rate variation in *Drosophila melanogaster*

Andrew H. Chan^{1,*}, Paul A. Jenkins^{1,*}, Yun S. Song^{1,2,**}

¹ Computer Science Division, University of California, Berkeley, CA, USA

² Department of Statistics, University of California, Berkeley, CA, USA

* These authors contributed equally to this work

** Corresponding author e-mail: yss@cs.berkeley.edu

Text S1

Two-locus recursion relation

Suppose we sample n haplotypes, observing their alleles at each of two loci and obtaining configuration $\mathbf{n} = (\mathbf{a}, \mathbf{b}, \mathbf{c})$. Here $\mathbf{c} = (c_{ij})$ is a matrix of the counts of haplotypes for which both alleles were observed; c_{ij} is the number of haplotypes with allele i at the first locus and allele j at the second locus. We also allow for the possibility that a haplotype had data missing at one locus: $\mathbf{a} = (a_i)_{i=1,\dots,K}$ is the vector of counts of haplotypes with allele i observed at the first locus and missing data at the second locus, and $\mathbf{b} = (b_j)_{j=1,\dots,L}$ is the vector of counts of haplotypes with allele j observed at the second locus and missing data at the first locus. Further, let:

$$\begin{aligned} a &= \sum_{i=1}^K a_i, & c_{i\cdot} &= \sum_{j=1}^L c_{ij}, & c &= \sum_{i=1}^K \sum_{j=1}^L c_{ij}, \\ b &= \sum_{j=1}^L b_j, & c_{\cdot j} &= \sum_{i=1}^K c_{ij}, & n &= a + b + c. \end{aligned}$$

The probability that, when we sample n haplotypes in some fixed order, we obtain a set consistent with configuration \mathbf{n} , is denoted by $q(\mathbf{n}; \theta_A, \theta_B, \rho)$. This probability is a function of θ_A , θ_B , and ρ : the mutation rates at the two loci, and the recombination rate between them. The respective mutation transition matrices at the two loci, which we denote \mathbf{P}^A and \mathbf{P}^B , are fixed. A system of equations for $q(\mathbf{n}; \theta_A, \theta_B, \rho)$ is given in [1]. We denote by $q(\mathbf{n}, s_1, s_2; \theta_A, \theta_B, \rho)$ the joint probability of obtaining \mathbf{n} with the events that there were precisely s_1 mutations in the history of the sample at the first locus and s_2 mutations in the history of the sample at the second locus. The corresponding system of equations for $q(\mathbf{n}, s_1, s_2; \theta_A, \theta_B, \rho)$ is:

$$\begin{aligned} &[n(n-1) + \theta_A(a+c) + \theta_B(b+c) + \rho c]q((\mathbf{a}, \mathbf{b}, \mathbf{c}), s_1, s_2; \theta_A, \theta_B, \rho) = \\ &\sum_{i=1}^K a_i(a_i - 1 + 2c_{i\cdot})q((\mathbf{a} - \mathbf{e}_i, \mathbf{b}, \mathbf{c}), s_1, s_2; \theta_A, \theta_B, \rho) + \sum_{j=1}^L b_j(b_j - 1 + 2c_{\cdot j})q((\mathbf{a}, \mathbf{b} - \mathbf{e}_j, \mathbf{c}), s_1, s_2; \theta_A, \theta_B, \rho) \\ &+ \sum_{i=1}^K \sum_{j=1}^L [c_{ij}(c_{ij} - 1)q((\mathbf{a}, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij}), s_1, s_2; \theta_A, \theta_B, \rho) + 2a_i b_j q((\mathbf{a} - \mathbf{e}_i, \mathbf{b} - \mathbf{e}_j, \mathbf{c} + \mathbf{e}_{ij}), s_1, s_2; \theta_A, \theta_B, \rho)] \\ &+ \theta_A \sum_{i=1}^K \left[\sum_{j=1}^L c_{ij} \sum_{t=1}^K P_{ti}^A q((\mathbf{a}, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij} + \mathbf{e}_{tj}), s_1 - 1, s_2; \theta_A, \theta_B, \rho) \right. \end{aligned}$$

$$\begin{aligned}
& + a_i \sum_{t=1}^K P_{ti}^A q((\mathbf{a} - \mathbf{e}_i + \mathbf{e}_t, \mathbf{b}, \mathbf{c}), s_1 - 1, s_2; \theta_A, \theta_B, \rho) \Big] \\
& + \theta_B \sum_{j=1}^L \left[\sum_{i=1}^K c_{ij} \sum_{t=1}^L P_{tj}^B q((\mathbf{a}, \mathbf{b}, \mathbf{c} - \mathbf{e}_{ij} + \mathbf{e}_{it}), s_1, s_2 - 1; \theta_A, \theta_B, \rho) \right. \\
& \quad \left. + b_j \sum_{t=1}^L P_{tj}^B q((\mathbf{a}, \mathbf{b} - \mathbf{e}_j + \mathbf{e}_t, \mathbf{c}), s_1, s_2 - 1; \theta_A, \theta_B, \rho) \right] \\
& + \rho \sum_{i=1}^K \sum_{j=1}^L c_{ij} q((\mathbf{a} + \mathbf{e}_i, \mathbf{b} + \mathbf{e}_j, \mathbf{c} - \mathbf{e}_{ij}), s_1, s_2; \theta_A, \theta_B, \rho), \tag{1}
\end{aligned}$$

where \mathbf{e}_{ij} is a unit matrix whose (i, j) th entry is one and the rest are zero. As before, we suppose that we know the identity of the ancestral allele at each locus, say λ_A and λ_B at locus A and B, respectively. Then we replace the relevant instances of (1) with the following:

$$\begin{aligned}
q((\mathbf{0}, \mathbf{b}, \mathbf{e}_{ij}), s_1, s_2; \theta_A, \theta_B, \rho) &= \begin{cases} q((\mathbf{0}, \mathbf{b} + \mathbf{e}_j, \mathbf{0}), 0, s_2; \theta_A, \theta_B, \rho) & \text{if } i = \lambda_A \text{ and } s_1 = 0, \\ 0 & \text{otherwise,} \end{cases} \\
q((\mathbf{a}, \mathbf{0}, \mathbf{e}_{ij}), s_1, s_2; \theta_A, \theta_B, \rho) &= \begin{cases} q((\mathbf{a} + \mathbf{e}_i, \mathbf{0}, \mathbf{0}), s_1, 0; \theta_A, \theta_B, \rho) & \text{if } j = \lambda_B \text{ and } s_2 = 0, \\ 0 & \text{otherwise,} \end{cases} \\
q((\mathbf{e}_i, \mathbf{0}, \mathbf{0}), s_1, s_2; \theta_A, \theta_B, \rho) &= \begin{cases} 1 & \text{if } i = \lambda_A \text{ and } s_1 = s_2 = 0, \\ 0 & \text{otherwise,} \end{cases} \\
q((\mathbf{0}, \mathbf{e}_j, \mathbf{0}), s_1, s_2; \theta_A, \theta_B, \rho) &= \begin{cases} 1 & \text{if } j = \lambda_B \text{ and } s_1 = s_2 = 0, \\ 0 & \text{otherwise.} \end{cases} \tag{2}
\end{aligned}$$

Padé summation

Modifications to the approach described in [2] are made, following from the boundary conditions given above. These can be converted into modifications of entries of the dynamic programming tables given in [2]. For example, using (2) we have that

$$\begin{aligned}
q((\mathbf{a}, \mathbf{0}, \mathbf{e}_{i\lambda_B}), 1, 0; \theta_A, \theta_B, \rho) &= q((\mathbf{a} + \mathbf{e}_i, \mathbf{0}, \mathbf{0}), 1, 0; \theta_A, \theta_B, \rho) \\
&= q(\mathbf{a} + \mathbf{e}_i, 1; \theta_A) + \frac{0}{\rho} + \frac{0}{\rho^2} + \dots,
\end{aligned}$$

where $q(\mathbf{a} + \mathbf{e}_i, 1; \theta_A)$ is the one-locus solution given by equation (3) in the main text. Notice that this expansion is in fact independent of ρ , from which it follows (by comparison with eq. (3.7) of [2]) that a number of entries in the dynamic programming tables are modified. For example, the second row in the dynamic programming table for the configuration $(\mathbf{a}, \mathbf{0}, \mathbf{e}_{i\lambda_B})$ is set to zero. Other boundary conditions may be interpreted in a similar fashion.

Ancestral allele estimation

Suppose we have one genomic sequence of *D. simulans* and n sequences of *D. melanogaster*. Let S represent the sequence of *D. simulans* and $M^{(k)}$ represent the sequence of the k th *D. melanogaster*, where S_l denotes the l th base of the sequence, and $S_{\hat{l}}$ represents the sequence with the exclusion of the l th base. Given $(S, M^{(k)})$, let $T_l^{(k)}$ be the time to the most recent common ancestor (TMRCA) at locus l ; $f_l^{(k)}(t | M_{\hat{l}}, S_{\hat{l}})$ be the density of the TMRCA conditioned on both their sequences but *excluding* the l th locus; and $A_l^{(k)}$ be the ancestral allele at the l th locus, i.e., the allele of the most recent common ancestor (MRCA).

To compute the distribution on the ancestral allele at the l th locus conditioned on $M^{(k)}$ and S , we use

Bayes' theorem to obtain

$$\begin{aligned}
& \mathbb{P}(A_l^{(k)} = i \mid M^{(k)}, S) \\
&= \frac{\int_0^\infty p(A_l^{(k)} = i, M^{(k)}, S, T_l^{(k)} = t) dt}{\mathbb{P}(M^{(k)}, S)} \\
&= \frac{\int_0^\infty \mathbb{P}(M_l^{(k)}, S_l^{(k)} \mid A_l^{(k)} = i, T_l^{(k)}) p(A_l^{(k)} = i, T_l^{(k)} = t) dt}{\mathbb{P}(M^{(k)}, S)} \\
&= \frac{\int_0^\infty \mathbb{P}(M_l^{(k)} \mid A_l^{(k)} = i, T_l^{(k)} = t) \mathbb{P}(S_l \mid A_l^{(k)} = i, T_l^{(k)} = t) \mathbb{P}(A_l^{(k)} = i) f_l^{(k)}(t \mid M_l^{(k)}, S_l) dt}{\sum_j \int_0^\infty \mathbb{P}(M_l^{(k)} \mid A_l^{(k)} = j, T_l^{(k)} = t) \mathbb{P}(S_l \mid A_l^{(k)} = j, T_l^{(k)} = t) \mathbb{P}(A_l^{(k)} = j) f_l^{(k)}(t \mid M_l^{(k)}, S_l) dt}. \quad (3)
\end{aligned}$$

In equation (3), the prior on the ancestral allele at locus l , $\mathbb{P}(A_l^{(k)} = i)$, is given by the stationary distribution of the allele frequencies from the mutation matrix \mathbf{P} . (In the above, p denotes a joint probability of discrete events together with the density for $T_l^{(k)}$.) The density on the TMRCA, $f_l^{(k)}(t \mid M_l^{(k)}, S_l)$, is estimated using Li & Durbin's `psmc` [3]. In practice, we use `psmc` to compute $f_l^{(k)}(t \mid M^{(k)}, S)$ and assume $f_l^{(k)}(t \mid M^{(k)}, S) \approx f_l^{(k)}(t \mid M_l^{(k)}, S_l)$.

The remaining two probabilities, $\mathbb{P}(M_l^{(k)} \mid A_l^{(k)} = i, T_l^{(k)} = t)$ and $\mathbb{P}(S_l \mid A_l^{(k)} = i, T_l^{(k)} = t)$, are computed as follows. For the computation of $\mathbb{P}(M_l^{(k)} \mid A_l^{(k)} = i, T_l^{(k)} = t)$, let $\mathbf{P} = (P_{ij})$ denote the mutation matrix, and let $r_l^{(k)}$ specify the number of mutations that have occurred at the l th locus of the k th *D. melanogaster* sequence during time $T_l^{(k)}$. Then we have

$$\begin{aligned}
\mathbb{P}(M_l^{(k)} = j \mid A_l^{(k)} = i, T_l^{(k)} = t) &= \sum_{s=0}^\infty \mathbb{P}(r_l^{(k)} = s \mid T_l^{(k)} = t) (\mathbf{P}^s)_{ij} \\
&= \sum_{s=0}^\infty \left(\frac{\theta t}{2} \right)^s \frac{e^{-\theta t/2}}{s!} (\mathbf{P}^s)_{ij} \\
&= \sum_{s=0}^\infty \left[\left(\frac{\theta t}{2} \mathbf{P} \right)^s \right]_{ij} \frac{e^{-\theta t/2}}{s!} \\
&= \left[e^{\frac{\theta t}{2} (\mathbf{P} - \mathbf{I})} \right]_{ij},
\end{aligned}$$

where \mathbf{I} is the identity matrix with the same dimensions as \mathbf{P} . The computation for $\mathbf{P}(S_l \mid A_l^{(k)} = j, T_l^{(k)} = t)$ is analogous.

After computing $\mathbb{P}(A_l^{(k)} = i \mid M^{(k)}, S)$ for every k and given l , we heuristically aggregate these pairwise probabilities to estimate $\mathbb{P}(A_l^{(k)} = i \mid M^{(1)}, \dots, M^{(n)}, S)$ as follows. Let $\bar{t}_l^{(k)}$ be the posterior mean of $f_l^{(k)}(t \mid M^{(k)}, S)$, i.e.:

$$\bar{t}_l^{(k)} = \int_0^\infty t f_l^{(k)}(t \mid M^{(k)}, S) dt,$$

and define $\tau_l = \max_k \bar{t}_l^{(k)}$. We approximate $\mathbb{P}(A_l^{(k)} = i \mid M^{(1)}, \dots, M^{(n)}, S)$ as

$$\mathbb{P}(A_l^{(k)} = i \mid M^{(1)}, \dots, M^{(n)}, S) \approx \frac{\sum_{k=1}^n \mathbb{P}(A_l^{(k)} = i \mid M^{(k)}, S) f_l^{(k)}(\tau_l \mid M_l^{(k)}, S_l)}{\sum_j \sum_{k=1}^n \mathbb{P}(A_l^{(k)} = j \mid M^{(k)}, S) f_l^{(k)}(\tau_l \mid M_l^{(k)}, S_l)},$$

which is a weighted average of $\mathbb{P}(A_l^{(k)} = i \mid M^{(k)}, S)$ over k , weighted by the density of the TMRCA evaluated at τ_l for each k . This averaging ought to mitigate effects such as genotyping errors and incomplete lineage sorting in individual *D. melanogaster* genomes.

References

1. Jenkins PA, Song YS (2009) Closed-form two-locus sampling distributions: accuracy and universality. *Genetics* 183: 1087–1103.
2. Jenkins PA, Song YS (2012) Padé approximants and exact two-locus sampling distributions. *Annals of Applied Probability* 22: 576–607.
3. Li H, Durbin R (2011) Inference of human population history from individual whole-genome sequences. *Nature* 475: 493–496.