

## 9.5 Estimation of pairwise coalescence times

Regev Schweiger and Richard Durbin (private communication) recently observed that the posterior mean of TMRCA in PSMC (Li and Durbin, 2011) is biased, and that the extent of bias depends on the true TMRCA as well as mutation and recombination rates. In what follows, we try to explain this phenomenon theoretically. For simplicity, we do not consider a full hidden Markov model, but the analysis carried out below captures the essence of the issue. Indeed, the theoretical predictions (Figure 9.7–Figure 9.9) derived below agree well with the empirical results that Schweiger and Durbin obtained using PSMC.

### 9.5.1 Preliminaries

Assume a constant-size, panmictic population and the infinite-sites model where mutations arrive according to a Poisson Point Process with rate  $\theta/2$ . Let  $T$  denote the TMRCA for a sample of size 2 and let  $M$  denote the number of mutations on the 2-leaved coalescent tree. Then, we have

$$\mathbb{P}[M = m \mid T = t] = \frac{1}{m!}(\theta t)^m e^{-\theta t}, \quad (9.5)$$

while the posterior density of  $T$  given  $M = m \in \mathbb{N}_0$  is (see Proposition 3.5)

$$f(t \mid M = m) = \frac{1}{m!}(1 + \theta)^{m+1} t^m e^{-(1+\theta)t}. \quad (9.6)$$

Hence, if  $\theta$  is known and we are given the observation  $M = m$ , then we can construct the following estimators of  $T$ :

1. **Maximum likelihood:** By maximizing (9.5), we get

$$\hat{T}_{\text{ML}} = \frac{m}{\theta}.$$

2. **Maximum a posteriori:** By maximizing (9.6), we get

$$\hat{T}_{\text{MAP}} = \frac{m}{1 + \theta}.$$

3. **Posterior mean:** By taking the expectation with respect to (9.6), we get

$$\hat{T}_{\text{PM}} = \frac{m + 1}{1 + \theta}.$$

Now, we note that the parameter  $\theta$  is in fact a random variable and that it depends on the size (or width)  $W$  of the genomic region that supports the coalescent tree with TMRCA  $T$ . More precisely,

$$\theta = \theta_0 W,$$

where  $\theta_0 = 4N_e\mu$  with  $\mu$  being the per-base per-generation mutation rate. The distribution of  $W$  depends on the TMRCA  $T$  and the recombination rate. Assume that the recombination rate is constant over the region of interest and define  $\rho = 4N_e r$ , where  $r$  corresponds to the per-generation recombination rate between two consecutive bases. Then,

$$\mathbb{P}[W = w \mid T = t] = [1 - p(\rho, t)]^{w-1} p(\rho, t),$$

where  $p(\rho, t)$  denotes the success probability of changing the TMRCA in one-base transition step. Under SMC' (Marjoram and Wall, 2006), a generalized version of the sequentially Markov coalescent,  $p(\rho, t)$  can be found as (Hobolth and Jensen, 2014)

$$p(\rho, t) = [e^{tQ}]_{1,1}, \quad \text{where } Q = \begin{pmatrix} -\rho & \rho & 0 \\ 1 & -2 & 1 \\ 0 & 0 & 0 \end{pmatrix}. \quad (9.7)$$

### 9.5.2 Bias

Putting everything together, the conditional expectation of an estimator  $\hat{T}$  given  $T = t$  can be found as

$$\begin{aligned} \mathbb{E}[\hat{T}(M, W) \mid T = t] &= \sum_{w=1}^{\infty} \sum_{m=0}^{\infty} \hat{T}(m, w) \mathbb{P}[M = m \mid T = t, W = w] \mathbb{P}[W = w \mid T = t] \\ &= \sum_{w=1}^{\infty} \sum_{m=0}^{\infty} \hat{T}(m, w) \left[ \frac{1}{m!} (\theta_0 w t)^m e^{-\theta_0 w t} \right] [1 - p(\rho, t)]^{w-1} p(\rho, t). \end{aligned} \quad (9.8)$$

For  $\hat{T} = \hat{T}_{\text{ML}}$ , the above computation can be carried out in closed-form and one obtains

$$\mathbb{E}[\hat{T}_{\text{ML}}(M, W) \mid T = t] = t. \quad (9.9)$$

For  $\hat{T} = \hat{T}_{\text{MAP}}$ , we obtain

$$\mathbb{E}[\hat{T}_{\text{MAP}}(M, W) \mid T = t] = \frac{\theta_0 t}{1 + \theta_0} {}_2F_1\left(1, \frac{1}{\theta_0}; 2 + \frac{1}{\theta_0}; 1 - p\right), \quad (9.10)$$

where  ${}_2F_1$  denotes the ordinary hypergeometric function. Finally, for  $\hat{T} = \hat{T}_{\text{PM}}$ , we get

$$\mathbb{E}[\hat{T}_{\text{PM}}(M, W) \mid T = t] = \frac{\theta_0 t}{1 + \theta_0} {}_2F_1\left(1, \frac{1}{\theta_0}; 2 + \frac{1}{\theta_0}; 1 - p\right) + \frac{p}{\theta_0} \Phi\left(1 - p, 1, 1 + \frac{1}{\theta_0}\right), \quad (9.11)$$

where  $\Phi$  denotes the Hurwitz-Lerch transcendent.

We see from (9.9) that  $\hat{T}_{\text{ML}}$  is an unbiased estimator of  $T$ , whereas (9.10) and (9.11) imply that  $\hat{T}_{\text{MAP}}$  and  $\hat{T}_{\text{PM}}$  are biased. Figures 9.7 and 9.8 respectively illustrate how  $\mathbb{E}[\hat{T}_{\text{MAP}}(M, W) \mid T = t]$  and  $\mathbb{E}[\hat{T}_{\text{PM}}(M, W) \mid T = t]$  deviate from  $t$ , for various values of  $r$  and  $\mu$ . In general,  $\hat{T}_{\text{PM}}$  appears less biased than  $\hat{T}_{\text{MAP}}$ .

### 9.5.3 Mean squared error

We saw in the previous section that the maximum likelihood estimator  $\hat{T}_{\text{ML}}$  is unbiased, while  $\hat{T}_{\text{MAP}}$  and  $\hat{T}_{\text{PM}}$  are biased. Does this imply that  $\hat{T}_{\text{ML}}$  is a better estimator than are

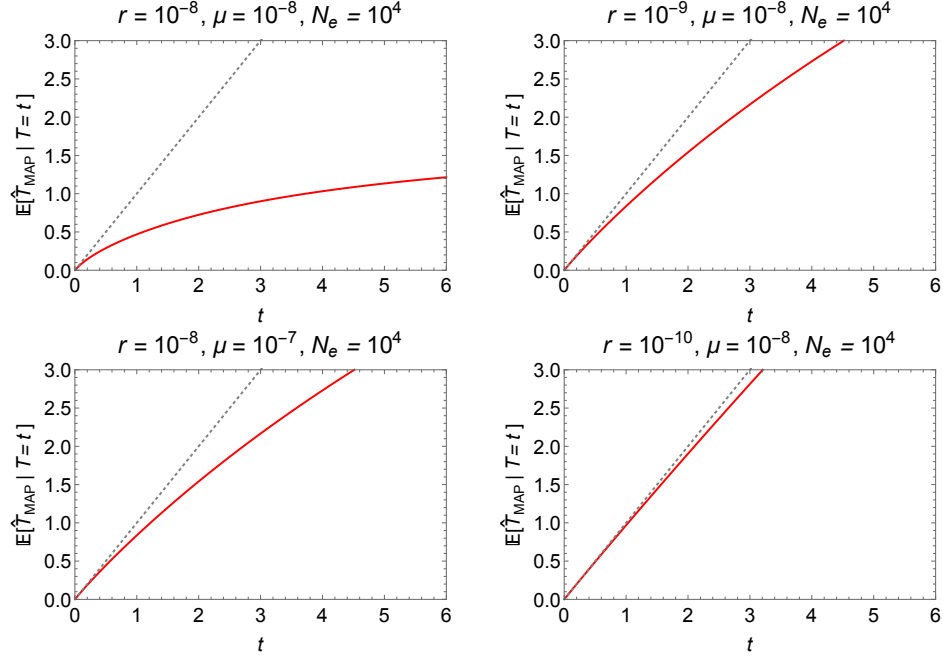


Fig. 9.7:  $\mathbb{E}[\hat{T}_{\text{MAP}}(M, W) \mid T = t]$  as a function of  $t$  for various values of  $r$  and  $\mu$ , with  $N_e = 10^4$ . Dotted lines correspond to  $y = x$ .

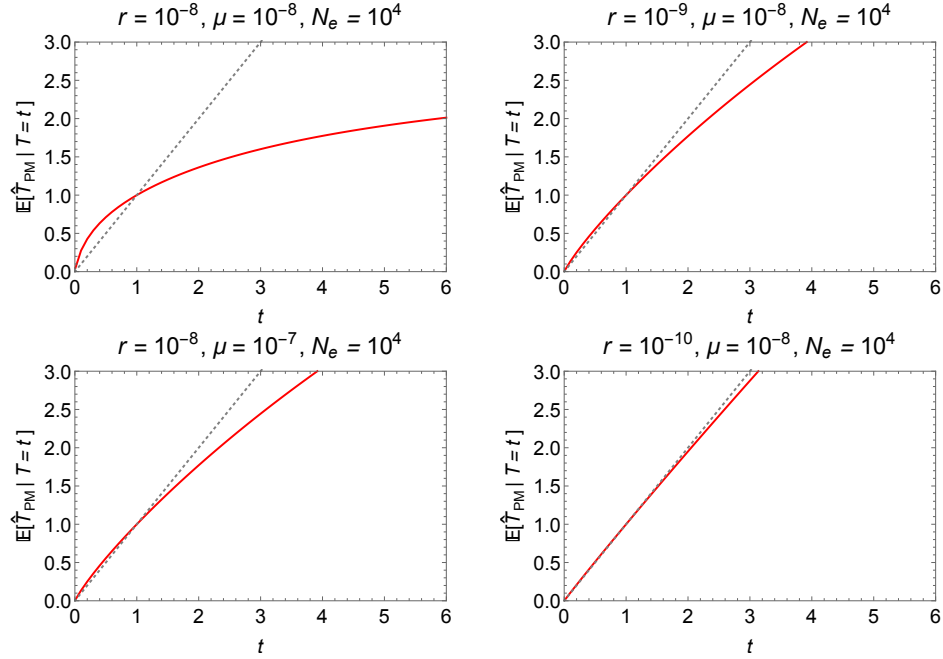


Fig. 9.8:  $\mathbb{E}[\hat{T}_{\text{PM}}(M, W) \mid T = t]$  as a function of  $t$  for various values of  $r$  and  $\mu$ , with  $N_e = 10^4$ . Dotted lines correspond to  $y = x$ .

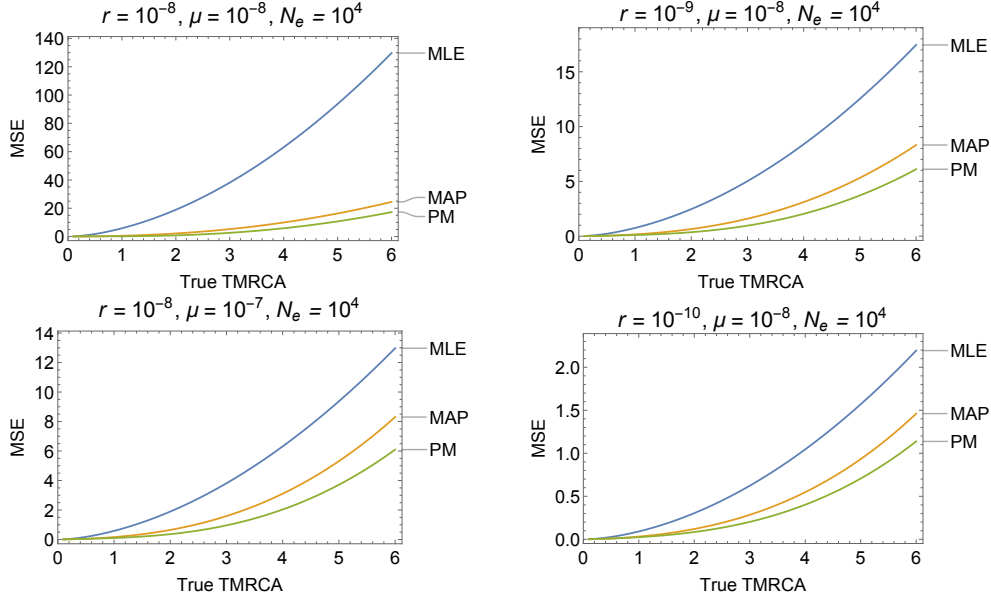


Fig. 9.9: Mean squared error (MSE) as a function of  $t$  for various values of  $r$  and  $\mu$ , with  $N_e = 10^4$ .

$\hat{T}_{\text{MAP}}$  and  $\hat{T}_{\text{PM}}$ ? Not necessarily! To see why, consider a quadratic loss function, in which case the risk  $R(t, \hat{T})$  corresponds to the mean squared error (MSE), which can be decomposed as

$$R(t, \hat{T}) = \text{MSE}[\hat{T}] = \mathbb{E}[(\hat{T} - t)^2] = \text{Var}(\hat{T}) + [\text{bias}(\hat{T})]^2,$$

where  $\text{bias}(\hat{T}) = \mathbb{E}[\hat{T}(M, W) \mid T = t] - t$ . So, the key question is how  $\text{Var}(\hat{T}_{\text{ML}})$  compares with  $\text{Var}(\hat{T}_{\text{MAP}})$  and  $\text{Var}(\hat{T}_{\text{PM}})$ . We can compute  $\mathbb{E}[(\hat{T} - t)^2 \mid T = t]$  in a similar vein as in (9.8), yielding

$$\begin{aligned} R(t, \hat{T}_{\text{ML}}) &= \frac{-tp \log(p)}{(1-p)\theta_0}, \\ R(t, \hat{T}_{\text{MAP}}) &= \frac{pt}{\theta_0^2} \left[ \theta_0 \Phi\left(1-p, 1, 1 + \frac{1}{\theta_0}\right) + (t-1)\Phi\left(1-p, 2, 1 + \frac{1}{\theta_0}\right) \right], \\ R(t, \hat{T}_{\text{PM}}) &= \frac{p}{\theta_0^2} \left[ \theta_0 t \Phi\left(1-p, 1, 1 + \frac{1}{\theta_0}\right) + [1 + t(t-3)]\Phi\left(1-p, 2, 1 + \frac{1}{\theta_0}\right) \right], \end{aligned}$$

where  $p$  depends on the population-scaled recombination rate  $\rho$  and the true TMRCA  $t$  as shown in (9.7). As Figure 9.9 shows, it turns out that  $\text{Var}(\hat{T}_{\text{ML}})$  can be substantially larger than  $\text{Var}(\hat{T}_{\text{MAP}})$  and  $\text{Var}(\hat{T}_{\text{PM}})$  if  $t \gg 0$ , thereby leading to much larger expected loss  $R(t, \hat{T}_{\text{ML}})$  compared to  $R(t, \hat{T}_{\text{MAP}})$  and  $R(t, \hat{T}_{\text{PM}})$ .

## References

- Griffiths RC (1981) Neutral two-locus multiple allele models with recombination. *Theoretical Population Biology* 19:169–186
- Griffiths RC (1991) The two-locus ancestral graph. *Selected Proceedings of the Sheffield Symposium on Applied Probability IMS Lecture Notes–Monograph Series* 18:100–117
- Griffiths RC, Marjoram P (1997) An ancestral recombination graph. In: Donnelly P, Tavaré S (eds) *Progress in population genetics and human evolution*, vol 87, Springer-Verlag, Berlin, pp 257–270
- Hobolth A, Jensen JL (2014) Markovian approximation to the finite loci coalescent with recombination along multiple sequences. *Theoretical population biology* 98:48–58
- Hudson RR (1983) Properties of a neutral allele model with intragenic recombination. *Theoretical Population Biology* 23:183–201
- Li H, Durbin R (2011) Inference of human population history from individual whole-genome sequences. *Nature* 475:493–496
- Marjoram P, Wall JD (2006) Fast “coalescent” simulation. *BMC Genet* 7:16
- Paigen K, Petkov P (2010) Mammalian recombination hot spots: properties, control and evolution. *Nature Reviews Genetics* 11(3):221–233
- Sasaki M, Lange J, Keeney S (2010) Genome destabilization by homologous recombination in the germ line. *Nature Reviews Molecular Cell Biology* 11(3):182–195
- Wiuf C (2000) A coalescence approach to gene conversion. *Theoretical Population Biology* 57:357–367