

# Understand graph-prior ramifications

For a given recurrent state, we assume relationship graphs are uniformly distributed *a priori*. This has various ramifications discussed below.

## Contents

<b>1</b>	<b>Posterior bounds that are induced by the graph prior</b>	<b>1</b>
1.1	How bounds are derived . . . . .	1
1.1.1	Single recurrence . . . . .	1
1.1.2	More than one recurrence . . . . .	2
<b>2</b>	<b>Bounds as indicators of data informativeness</b>	<b>4</b>
<b>3</b>	<b>How bounds induced by the graph prior change with MOI</b>	<b>4</b>
<b>4</b>	<b>A counter-intuitive property</b>	<b>4</b>

## 1 Posterior bounds that are induced by the graph prior

When the prior probability on relapse is non-zero, posterior probabilities of reinfection and recrudescence never reach one because all genetic data is always compatible with relapse. In this section, we discuss bounds on posterior probabilities of reinfection and recrudescence; they are induced by the uniformity of the prior over relationship graphs. Bounds define a feasible set of posterior probabilities without any information about the genetic data beyond that used to derive the multiplicities of infection (MOIs). For example, in the case of two genotypes in an enrolment episode followed by a monoclonal recurrence, when recurrent states are equally likely *a priori*, the posterior probabilities of reinfection and recrudescence are upper bounded by  $9/11$  and  $9/13$ , respectively. The feasible set of posterior probabilities is shown in Figure 1.

### 1.1 How bounds are derived

We start by making some observations about the posterior odds of relapse to reinfection / recrudescence; odds help simplify some derivations.

#### 1.1.1 Single recurrence

For clarity of exposition, we start with the case of a single recurrent episode. Let  $\mathcal{G}_C, \mathcal{G}_L$ , and  $\mathcal{G}_I$  denote subsets of the graph space  $\mathcal{G}$ , containing the relationship graphs compatible with recrudescence, relapse, and reinfection respectively. (To be clear,  $\mathcal{G}_L$  is exactly  $\mathcal{G}$ .) The posterior odds of relapse to recrudescence is given by

$$o_{L:C} := \frac{\mathbb{P}(\mathbf{y}|L)\mathbb{P}(L)}{\mathbb{P}(\mathbf{y}|C)\mathbb{P}(C)} = \frac{\mathbb{P}(L)}{\mathbb{P}(C)} \frac{\sum_{g \in \mathcal{G}_L} \mathbb{P}(\mathbf{y}|g)\mathbb{P}(g|L)}{\sum_{g \in \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g)\mathbb{P}(g|C)} = \frac{\mathbb{P}(L)}{\mathbb{P}(C)} \frac{|\mathcal{G}_C|}{|\mathcal{G}_L|} \frac{\sum_{g \in \mathcal{G}_L} \mathbb{P}(\mathbf{y}|g)}{\sum_{g \in \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g)} = \frac{\mathbb{P}(L)}{\mathbb{P}(C)} \frac{|\mathcal{G}_C|}{|\mathcal{G}_L|} \left( 1 + \frac{\sum_{g \in \mathcal{G}_L \setminus \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g)}{\sum_{g \in \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g)} \right),$$

where  $\mathcal{G}_L \setminus \mathcal{G}_C$  is the subset of graphs compatible with relapse but not recrudescence (graphs where genotypes from the recurrent episode are not all clones from the initial episode). It follows that  $o_{L:C} \geq \frac{\mathbb{P}(L)|\mathcal{G}_C|/\mathbb{P}(C)|\mathcal{G}_L|}{1}$ . Note that this inequality is close to equality when  $\sum_{g \in \mathcal{G}_L \setminus \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g) \ll \sum_{g \in \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g)$ .

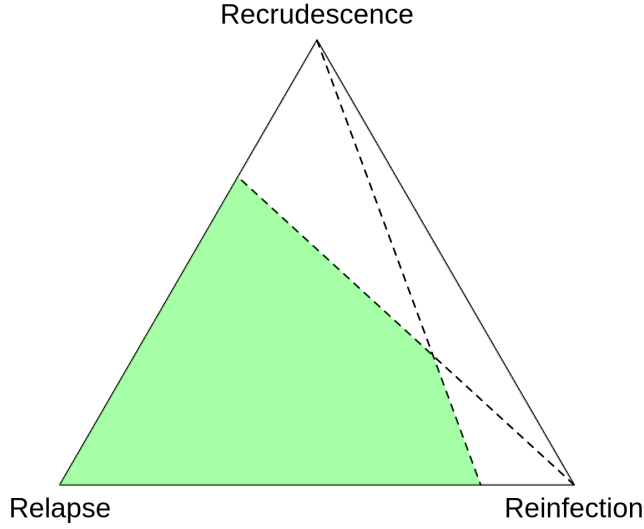


Figure 1: Feasible set of posterior probabilities in the case of two genotypes in the initial episode and one genotype in the recurrent episode. Dashed lines intersect the left and bottom edges of the simplex at  $\frac{4}{13}$  and  $\frac{2}{11}$ , respectively.

Similarly, the posterior odds of relapse to reinfection is given by

$$o_{L:I} := \frac{\mathbb{P}(\mathbf{y}|L)\mathbb{P}(L)}{\mathbb{P}(\mathbf{y}|I)\mathbb{P}(I)} = \frac{\mathbb{P}(L)}{\mathbb{P}(I)} \frac{|\mathcal{G}_I|}{|\mathcal{G}_L|} \left( 1 + \frac{\sum_{g \in \mathcal{G}_L \setminus \mathcal{G}_I} \mathbb{P}(\mathbf{y}|g)}{\sum_{g \in \mathcal{G}_I} \mathbb{P}(\mathbf{y}|g)} \right).$$

It follows that  $o_{L:I} \geq \mathbb{P}(L)|\mathcal{G}_I|/\mathbb{P}(I)|\mathcal{G}_L|$ . Note that this inequality is close to equality when  $\sum_{g \in \mathcal{G}_L \setminus \mathcal{G}_I} \mathbb{P}(\mathbf{y}|g) \ll \sum_{g \in \mathcal{G}_I} \mathbb{P}(\mathbf{y}|g)$ .

The posterior odds  $o_{I:C}$  can be defined similarly, but analogous lower bounds are not available. This is because  $\mathcal{G}_C$  and  $\mathcal{G}_I$  do not intersect. Thus, we simply note that  $o_{I:C} \approx 0$  when  $\mathbb{P}(\mathbf{y}|I) \ll \mathbb{P}(\mathbf{y}|C)$ .

As for the posterior probabilities, we have

$$\begin{aligned} \mathbb{P}(C|\mathbf{y}) &= \frac{\mathbb{P}(\mathbf{y}|C)\mathbb{P}(C)}{\mathbb{P}(\mathbf{y}|C)\mathbb{P}(C) + \mathbb{P}(\mathbf{y}|L)\mathbb{P}(L) + \mathbb{P}(\mathbf{y}|I)\mathbb{P}(I)} \\ &= \frac{1}{1 + o_{L:C} + o_{I:C}} \\ &\leq \frac{1}{1 + \frac{\mathbb{P}(L)|\mathcal{G}_C|}{\mathbb{P}(C)|\mathcal{G}_L|} + 0} \\ &= \frac{\mathbb{P}(C)}{\mathbb{P}(C) + \mathbb{P}(L) \frac{|\mathcal{G}_C|}{|\mathcal{G}_L|}}. \end{aligned} \tag{1}$$

This inequality is close to equality when  $\sum_{g \in \mathcal{G}_L \setminus \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g) \ll \sum_{g \in \mathcal{G}_C} \mathbb{P}(\mathbf{y}|g)$  and  $\mathbb{P}(\mathbf{y}|I) \ll \mathbb{P}(\mathbf{y}|C)$ . Note that the former implies the latter as  $\mathcal{G}_I \subseteq \mathcal{G}_L \setminus \mathcal{G}_C$ .

Similarly, we have

$$\mathbb{P}(I|\mathbf{y}) \leq \frac{\mathbb{P}(I)}{\mathbb{P}(I) + \mathbb{P}(L) \frac{|\mathcal{G}_I|}{|\mathcal{G}_L|}}. \tag{2}$$

Likewise, this inequality is close to equality when  $\sum_{g \in \mathcal{G}_L \setminus \mathcal{G}_I} \mathbb{P}(\mathbf{y}|g) \ll \sum_{g \in \mathcal{G}_I} \mathbb{P}(\mathbf{y}|g)$ .

### 1.1.2 More than one recurrence

When there is more than one recurrence and prior probabilities are non-zero, the maximum probabilities of all sequences are non-certain with the exception of the relapse-only sequence. As before, we can compute odds; for

example

$$o_{LL:CL} := \frac{\mathbb{P}(\mathbf{y}|LL)\mathbb{P}(LL)}{\mathbb{P}(\mathbf{y}|CL)\mathbb{P}(CL)} = \frac{\mathbb{P}(LL)}{\mathbb{P}(CL)} \frac{|\mathcal{G}_{CL}|}{|\mathcal{G}_{LL}|} \frac{\sum_{\mathbf{g} \in \mathcal{G}_{LL}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{\sum_{\mathbf{g} \in \mathcal{G}_{CL}} \mathbb{P}(\mathbf{y}|\mathbf{g})} = \frac{\mathbb{P}(LL)}{\mathbb{P}(CL)} \frac{|\mathcal{G}_{CL}|}{|\mathcal{G}_{LL}|} \left( 1 + \frac{\sum_{\mathbf{g} \in \mathcal{G}_{LL} \setminus \mathcal{G}_{CL}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{\sum_{\mathbf{g} \in \mathcal{G}_{CL}} \mathbb{P}(\mathbf{y}|\mathbf{g})} \right).$$

However, unlike before, we cannot always derive maximum posterior probabilities from the odds because the probabilities of the remaining sequences are not necessarily zero (the exception being when the odds of an all-but-one-relapse to all-relapse sequence are maximised). Instead, we must compute maximum probabilities the long way; for example,

$$\begin{aligned} \mathbb{P}(CC|\mathbf{y}) = & \left\{ \sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(CC)/|\mathcal{G}_{CC}|} \right\} \left\{ \sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(CC)/|\mathcal{G}_{CC}|} + \left( \sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g}) + \sum_{\mathbf{g} \in \mathcal{G}_{LC} \setminus \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g}) \right) \frac{\mathbb{P}(LC)}{|\mathcal{G}_{LC}|} \right. \\ & + \left( \sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g}) + \sum_{\mathbf{g} \in \mathcal{G}_{LL} \setminus \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g}) \right) \frac{\mathbb{P}(LL)}{|\mathcal{G}_{LL}|} + \left( \sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g}) + \sum_{\mathbf{g} \in \mathcal{G}_{CL} \setminus \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g}) \right) \frac{\mathbb{P}(CL)}{|\mathcal{G}_{CL}|} + \sum_{\mathbf{g} \in \mathcal{G}_{CI}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(CI)/|\mathcal{G}_{CI}|} \\ & \left. + \sum_{\mathbf{g} \in \mathcal{G}_{LI}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(LI)/|\mathcal{G}_{LI}|} + \sum_{\mathbf{g} \in \mathcal{G}_{IC}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(IC)/|\mathcal{G}_{IC}|} + \sum_{\mathbf{g} \in \mathcal{G}_{IL}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(IL)/|\mathcal{G}_{IL}|} + \sum_{\mathbf{g} \in \mathcal{G}_{II}} \mathbb{P}(\mathbf{y}|\mathbf{g})^{\mathbb{P}(II)/|\mathcal{G}_{II}|} \right\}^{-1} \quad (3) \end{aligned}$$

is approximately

$$\frac{\mathbb{P}(CC)}{\mathbb{P}(CC) + \mathbb{P}(LC) \frac{|\mathcal{G}_{CC}|}{|\mathcal{G}_{LC}|} + \mathbb{P}(LL) \frac{|\mathcal{G}_{CC}|}{|\mathcal{G}_{LL}|} + \mathbb{P}(CL) \frac{|\mathcal{G}_{CC}|}{|\mathcal{G}_{CL}|}}$$

when  $\sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g})$  significantly exceeds the summation over all other subsets of graph space in equation (3).

Bounds on marginal posteriors are harder to compute because computation involves summation over sequences whose probabilities are not necessarily zero. For example, the bound on recrudescence at the first recurrence of two recurrences, may be close to one of three upper bounds each involving non-zero probabilities:

1. bound given maximal CC probability (probabilities of CI & CL are not necessarily both zero)
2. bound given maximal CL probability (probabilities of CC & CI are not necessarily both zero)
3. bound given maximal CI probability (probabilities of CC & CL are not necessarily both zero)

**Digression** In the online article “Understand posterior estimates” we compute maximum marginal probabilities of recrudescence / reinfection at the first recurrence of two / three recurrences. These maxima are not bounds imposed by the prior: they are based on additional knowledge that there are no recurrent data on all but the first recurrence. The computation only holds because all episodes are monoclonal. When all episodes are monoclonal, the likelihood of equivalent sequences (e.g., CC, CI and CL in the case of two recurrences with strong evidence of recrudescence on the first and no data on the second) are equal (otherwise, they are unequal for reasons analogous to those that explain departure from the prior in the section “Data on only one episode”) and graph likelihoods are equal for all graphs compatible with equivalent sequences (e.g.,  $\mathbb{P}(\mathbf{y}|\mathbf{g})$  is the same  $\forall \mathbf{g} \in \mathcal{G}_{CC}$ ,  $\forall \mathbf{g} \in \mathcal{G}_{CI}$ , and  $\forall \mathbf{g} \in \mathcal{G}_{CL}$  in the case of two recurrences with strong evidence of recrudescence on the first and no data on the second). Given a uniform prior on recurrent states, three monoclonal episodes with strong evidence of recrudescence for the first recurrence and no data on the second recurrence, we thus have

$$\mathbb{P}(s_1 = C|\mathbf{y}) \quad (4)$$

$$\begin{aligned} & \leq \frac{\sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g})/|\mathcal{G}_{CC}|}{\frac{\sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{|\mathcal{G}_{CC}|} + \frac{\sum_{\mathbf{g} \in \mathcal{G}_{CI}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{|\mathcal{G}_{CI}|} + \frac{\sum_{\mathbf{g} \in \mathcal{G}_{CL}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{|\mathcal{G}_{CL}|} + \frac{\sum_{\mathbf{g} \in \mathcal{G}_{CC}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{|\mathcal{G}_{LC}|} + \frac{\sum_{\mathbf{g} \in \mathcal{G}_{CI}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{|\mathcal{G}_{LI}|} + \frac{\sum_{\mathbf{g} \in \mathcal{G}_{CL}} \mathbb{P}(\mathbf{y}|\mathbf{g})}{|\mathcal{G}_{LL}|}} \\ & = \frac{|\mathcal{G}_{CC}|/|\mathcal{G}_{CC}|}{\frac{|\mathcal{G}_{CC}|}{|\mathcal{G}_{CC}|} + \frac{|\mathcal{G}_{CI}|}{|\mathcal{G}_{CI}|} + \frac{|\mathcal{G}_{CL}|}{|\mathcal{G}_{CL}|} + \frac{|\mathcal{G}_{CC}|}{|\mathcal{G}_{LC}|} + \frac{|\mathcal{G}_{CI}|}{|\mathcal{G}_{LI}|} + \frac{|\mathcal{G}_{CL}|}{|\mathcal{G}_{LL}|}} \\ & = \frac{1}{3 + 1/3 + 1/3 + 3/12}. \quad (5) \end{aligned}$$

The same reasoning applies given strong evidence of reinfection.

## 2 Bounds as indicators of data informativeness

By comparing probability estimates to their respective bounds, we might like to answer the question: could the posterior probability of my most probable state sequence be higher if we had data on more markers? However, we cannot guarantee convergence onto bounds. In the case of a marginal probability, the bound is the maximum of multiple upper bounds (see example above), and thus convergence is not guaranteed; in the case of a non-marginal probability, we suspect convergence might depend on the relationship graph from which the data are generated.

## 3 How bounds induced by the graph prior change with MOI

In this section we explore how maximum probabilities of recrudescence / reinfection given a single recurrence change with graph size, where maximum probabilities are those induced by the the graph prior. Maximum probabilities of reinfection increase with the size of the graph, both when MOIs of the enrolment episode and the first recurrence are equal (plot 2a, centre) and not (plot 2a, off-centre). Maximum probabilities of recrudescence increase with graph size when the MOIs of the enrolment episode and the first recurrence are equal (plot 2b, centre); they decrease with increasing disparity between the higher MOI of the enrolment episode and the lower MOI of the first recurrence (plot 2b, right of centre). Note that, when the MOI of first recurrence exceeds that of the enrolment episode, recrudescence has zero posterior probability because we assume under the Pv3Rs model that all parasites are detected and that there are no genotyping errors; as such, a recrudescence can be at most as diverse as the preceding episode.

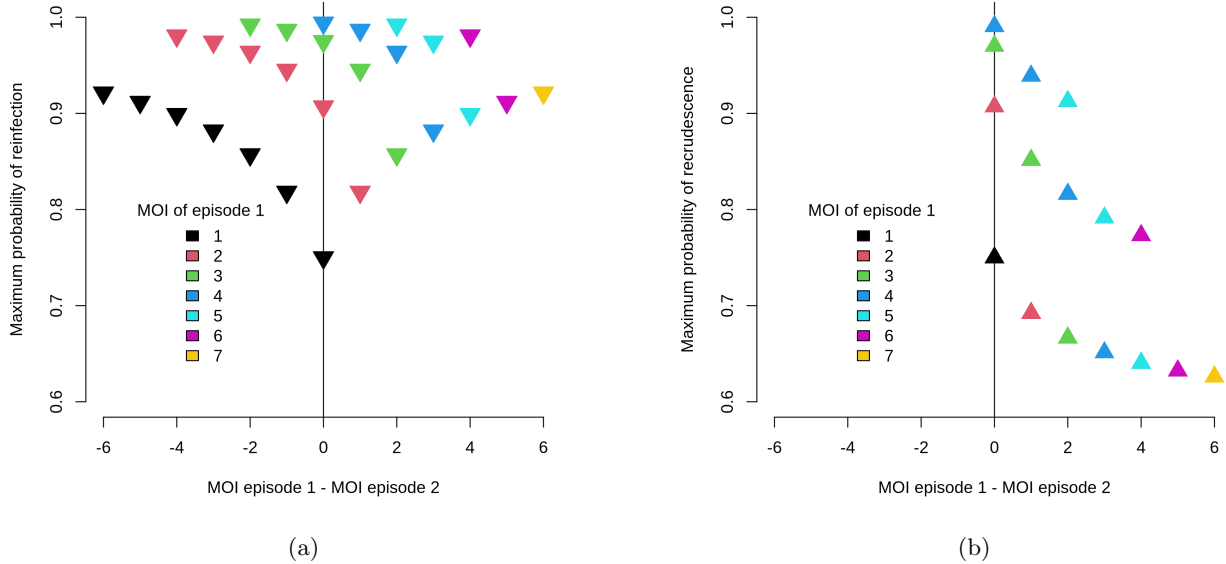


Figure 2: Maximum probabilities of recrudescence / reinfection given a single recurrence.

## 4 A counter-intuitive property

The *a priori* assumption that all valid relationships graphs  $\mathbf{g}$  are equally likely given a recurrent state leads to a counter-intuitive property. Specifically, the probability of an edge can vary depending on the graph it is embedded within. Consider the scenario where one monoclonal recurrence follows a monoclonal enrolment episode.

Because we assume that the three relationship graphs are equally likely given relapse, the probability distribution of the edge between the two genotypes given relapse is

$$\mathbb{P}(\text{the two genotypes are clones} | \mathbf{s} = L) = 1/3, \quad (6)$$

$$\mathbb{P}(\text{the two genotypes are siblings} | \mathbf{s} = L) = 1/3, \quad (7)$$

$$\mathbb{P}(\text{the two genotypes are strangers} | \mathbf{s} = L) = 1/3. \quad (8)$$

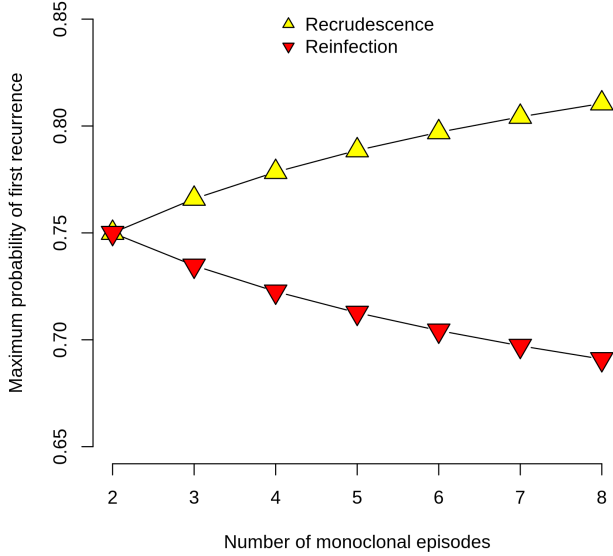


Figure 3: The probability that the first recurrence is a recrudescence / reinfection increases / decreases with additional recurrences devoid of data. The prior on recurrent states is uniform. The maximum probabilities are those based on knowledge that there are no recurrent data for all but the first recurrence and all episodes are monoclonal. For example, from equations (6) and (9) we have  $|\mathcal{G}_C|/|\mathcal{G}_L| = 1/3 > |\mathcal{G}_{CL}|/|\mathcal{G}_{LL}| = 3/12$ . Plugging  $|\mathcal{G}_C|/|\mathcal{G}_L| = 1/3$  into equation (1), and  $|\mathcal{G}_{CL}|/|\mathcal{G}_{LL}| = 3/12$  plus  $|\mathcal{G}_{CC}|/|\mathcal{G}_{LC}| = 1/3$  and  $|\mathcal{G}_{CI}|/|\mathcal{G}_{LI}| = 1/3$  into equation (5), we see that the probability that the first recurrence is a recrudescence when recurrent states are equally likely *a priori* increases slightly from  $3/4$  to  $3/(3 + 11/12)$  with the addition of the second recurrence without data.

Now consider the scenario in which there is an additional monoclonal recurrence.

There are now 12 relationship graphs, which are assumed to be equally likely under relapses. Among them, between the first two genotypes, three have a clonal edge, four have a sibling edge, and five have a stranger edge. The probability distribution of the edge between the first two genotypes given relapses is thus

$$\mathbb{P}(\text{the first two genotypes are clones} | \mathbf{s} = \text{LL}) = 3/12, \quad (9)$$

$$\mathbb{P}(\text{the first two genotypes are siblings} | \mathbf{s} = \text{LL}) = 4/12, \quad (10)$$

$$\mathbb{P}(\text{the first two genotypes are strangers} | \mathbf{s} = \text{LL}) = 5/12. \quad (11)$$

An explanation for the change in the probability distribution of the edge between the first two genotypes upon the addition of the second relapse is that a clonal edge between the first two genotypes imposes a constraint where the two remaining edges must exhibit the same relationship. Similarly, a sibling edge between the first two genotypes is incompatible with a single stranger edge among the two remaining edges. In general, the pattern that stranger edges are more likely *a priori* and clonal edges are less likely *a priori* becomes more prominent when more genotypes are present. For monoclonal episodes, this change in distribution over the edge between the first two genotypes results in an increase / decrease in the maximum probability that the first recurrence is a recrudescence / reinfection with the addition of monoclonal recurrences devoid of data (Figure 3), where the maximum probabilities are those assuming all episodes are monoclonal and there are no recurrent data on all but the first recurrence (e.g., as in equation (5)).