

Inferring population histories in the IM-model

Lars Nørvang Andersen Thomas Mailund Asger Hobolth

BiRC
Aarhus University

Montpellier, June 2012



Scally et al. *Nature*, 483 169–175, Mar 2012.

Gor

The case of 3 lineages

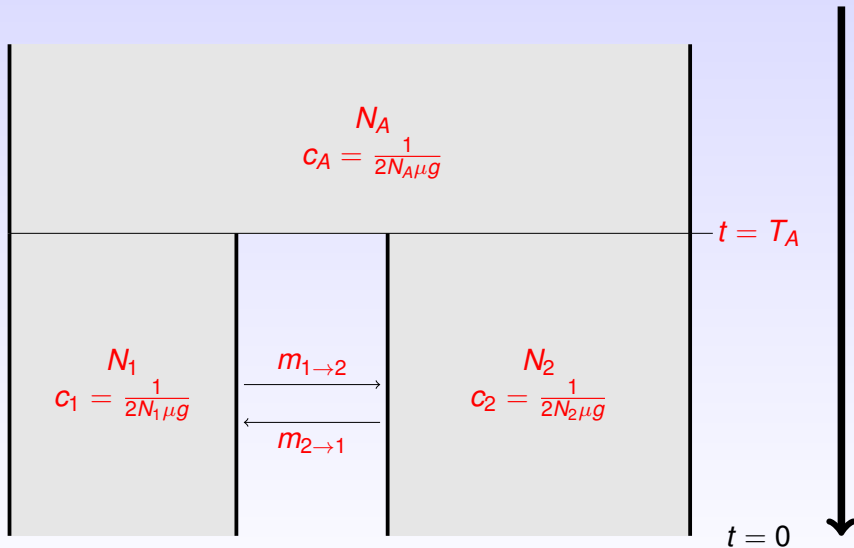
Our data points are the alignment columns:

A	A	A	...	T	T
A	A	A	...	T	T
A	C	G	...	C	T
N_{AAA}	N_{AAC}	N_{AAG}	...	N_{TTC}	N_{TTT}

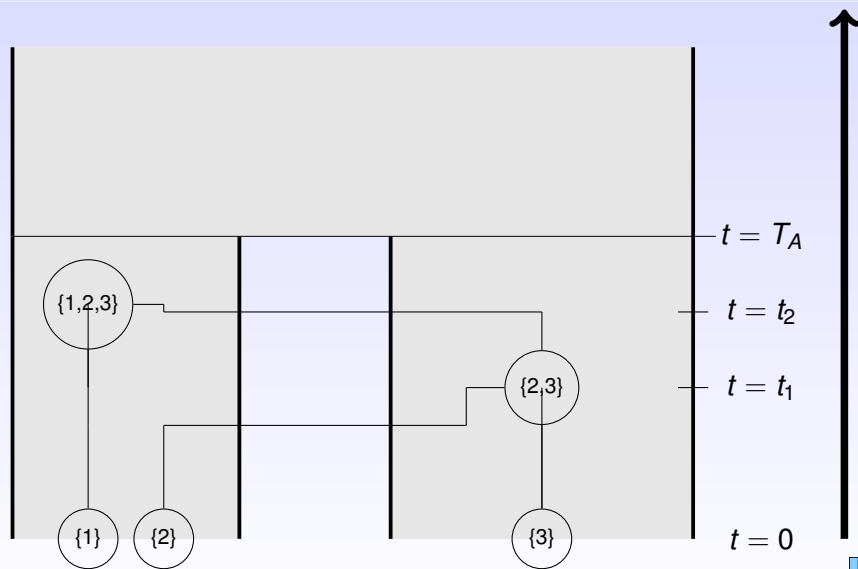
We want to analyze this data using maximum likelihood
Free recombination between loci implies that the log-likelihood function given the data \mathbf{x} is given by

$$L(\Theta | \mathbf{x}) = \sum_{x \in \{A, C, T, G\}^3} N(x) \log \mathbb{P}(x)$$

The IM-model



The IM-model



The IM-model - prior to T_A

The lineages are subsets $\ell_i \subseteq \{1, 2, \dots, L\}$

The populations are number $j_i \in \{1, 2, \dots, N\}$

Each state is a union of tuples: $\{(j_i, \ell_i) \mid i = 1, \dots, m\}$

Examples of states:

Starting state: $s = \{(1, \{1\}), (1, \{2\}), (2, \{3\})\}$

Migration: $s_1 = \{(1, \{1\}), (2, \{2\}), (2, \{3\})\}$

Coalescence: $s_2 = \{(1, \{1\}), (2, \{2, 3\})\}$

Migration: $s_3 = \{(1, \{1\}), (1, \{2, 3\})\}$

Coalescence: $s_4 = \{(1, \{1, 2, 3\})\}$

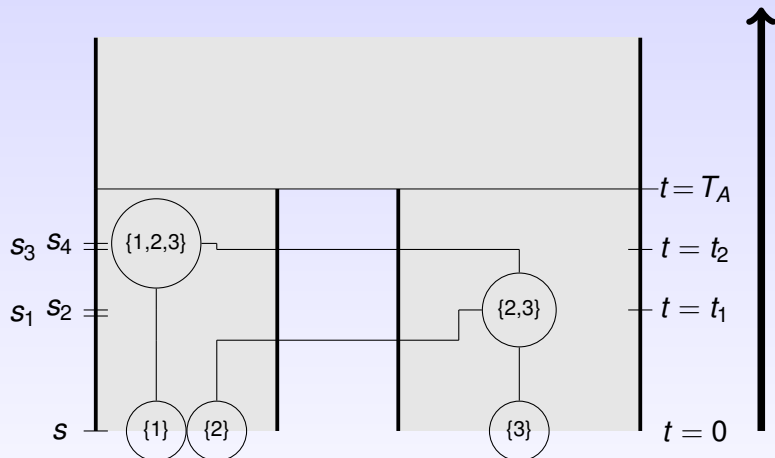


The IM-model - prior to T_A

$$Q_{\alpha,\beta} = \begin{cases} m_{i \rightarrow j} & \text{if } \alpha = \mathcal{S} \cup \{(i, l)\}, \beta = \mathcal{S} \cup \{(j, l)\} \\ c_i & \text{if } \alpha = \mathcal{S} \cup \{(i, l_1)\} \cup \{(i, l_2)\}, \beta = \mathcal{S} \cup \{(i, l_1 \cup l_2)\} \\ 0 & \text{otherwise,} \end{cases}$$

where \mathcal{S} is of the form $\cup_i \{(j_i, l_i)\}$

The IM-model

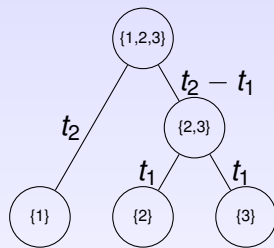


$s = \{(1, \{1\}), (1, \{2\}), (2, \{3\})\}$, $s_1 = \{(1, \{1\}), (2, \{2\}), (2, \{3\})\}$
 $s_2 = \{(1, \{1\}), (2, \{2,3\})\}$, $s_3 = \{(1, \{1\}), (1, \{2,3\})\}$, $s_4 = \{(1, \{1,2,3\})\}$

$$\left(e^{Qt_1} \right)_{s,s_1} c_2 \left(e^{Q(t_2-t_1)} \right)_{s_2,s_3} c_1$$

Genealogies - Coalescent trees

Sample paths of $X(t)$ are *genealogies*.
coalescent trees are unions of genealogies



$$f(C(t_1, t_2)) = \sum_{(\alpha_j)} \left(e^{Q t_1} \right)_{s, s_1} c_1 \quad \left(e^{Q(t_2 - t_1)} \right)_{s_2, s_3} c_2$$

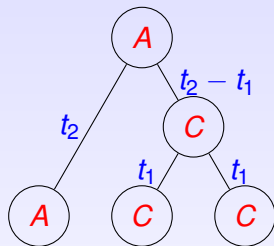
$$f(C(t_1, t_2)) = \sum_{(\alpha_j)} \left(e^{Q t_1} \right)_{s, \alpha_1} c_{(\alpha_1, \alpha_2)} \left(e^{Q(t_2 - t_1)} \right)_{\alpha_2, \alpha_3} c_{(\alpha_3, \alpha_4)}$$

The nucleotide substitution matrix - time-reversible case

$$\bar{Q} = \begin{array}{c} \\ A \\ C \\ G \\ T \end{array} \begin{array}{c} A \\ C \\ G \\ T \end{array} \begin{pmatrix} \cdot & \alpha\pi_C & \beta\pi_G & \gamma\pi_T \\ \alpha\pi_A & \cdot & \rho\pi_G & \sigma\pi_T \\ \beta\pi_A & \rho\pi_C & \cdot & \tau\pi_T \\ \gamma\pi_A & \sigma\pi_C & \tau\pi_G & \cdot \end{pmatrix}$$

Stationary distribution: $\pi = (\pi_A, \pi_C, \pi_G, \pi_T)$

Nucleotide assignment



$$\mathbb{P}(x \mid C(t_1, t_2)) = \pi_A \left(e^{\bar{Q}t_2} \right)_{A,A} \left(e^{\bar{Q}(t_2-t_1)} \right)_{A,C} \left(e^{\bar{Q}t_1} \right)_{C,C} \left(e^{\bar{Q}t_1} \right)_{C,C}$$

Explicit calculations

Recall:

$$f(C(t_1, t_2)) = \sum_{(\alpha_j)} \left(e^{Q t_1} \right)_{S, \alpha_1} C_{(\alpha_1, \alpha_2)} \left(e^{Q(t_2 - t_1)} \right)_{\alpha_2, \alpha_3} C_{(\alpha_3, \alpha_4)}$$

$$\begin{aligned} \mathbb{P}(x) &= \int \int \mathbb{P}(x \mid C(t_1, t_2)) f(C(t_1, t_2)) dt_1 dt_2 \\ &= \sum C \int \int \left(e^{\bar{Q} t_2} \right)_{A, A} \dots \left(e^{Q(t_2 - t_1)} \right)_{\alpha_2, \alpha_3} dt_1 dt_2 \\ &= \sum \sum C C_1 \int \int e^{\bar{\lambda}_i t_2} \dots e^{\lambda_j(t_2 - t_1)} dt_1 dt_2 \end{aligned}$$

$$\left(e^{\bar{Q} t} \right)_{x, y} = \sum_i \bar{v}_{x, i} \bar{v}_{i, y}^{-1} e^{\bar{\lambda}_i t},$$

Statespace explosion

The number of states increases very rapidly:

L	2	3	4	5	6
States	6	22	94	454	2430

Table : The number of states for $L = 2, 3, 4, 5, 6$ lineages

Given CTMC $\{X_t\}_{t \geq 0}, (S, Q, \pi_0)$, S state space, Q rate matrix, π_0 initial distribution.

Partition $\mathcal{P} = \{S_i\}_{i \in I}$

Define a stochastic process $\hat{X}_n = i \leftrightarrow X_n \in S_i$.

Under what conditions $\{\hat{X}_n\}$ a CTMC?

Strong Lumpability - Exact Lumpability

$$\mathcal{P} = \{\{S_1\}, \{S_2\}, \dots, \{S_N\}\}$$

$$\begin{pmatrix} (Q_{1,1}) & (Q_{1,2}) & \dots & (Q_{1,N}) \\ (Q_{2,1}) & (Q_{2,2}) & \dots & (Q_{2,N}) \\ \vdots & \vdots & \ddots & \vdots \\ (Q_{N,1}) & (Q_{N,2}) & \dots & (Q_{N,N}) \end{pmatrix} \begin{pmatrix} q_{1,1} & q_{1,2} & q_{1,3} & \dots & q_{1,n} \\ q_{2,1} & q_{2,2} & q_{2,3} & \dots & q_{2,n} \\ q_{3,1} & q_{3,2} & q_{3,3} & \dots & q_{3,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ q_{m,1} & q_{m,2} & q_{m,3} & \dots & q_{m,n} \end{pmatrix}$$

- Strong Lumpability : $\forall I \neq J, i, i' \in S_I : \sum_{j \in J} q_{i,j} = \sum_{j \in J} q_{i',j}$
- Exact Lumpability :
 $\forall S_I, S_J, j, j' \in S_J : \sum_{i \in S_I} q_{i,j} = \sum_{i \in S_I} q_{i,j'}$

Applying lumpability

Q :

$$\begin{array}{c}
 1\ 2\ 3 \mid \\
 1\ 2 \mid 3 \\
 1\ 3 \mid 2 \\
 2\ 3 \mid 1 \\
 1 \mid 2\ 3 \\
 2 \mid 1\ 3 \\
 3 \mid 1\ 2 \\
 \mid 1\ 2\ 3 \\
 A
 \end{array}
 \begin{pmatrix}
 \cdot & m_{1 \rightarrow 2} & m_{1 \rightarrow 2} & m_{1 \rightarrow 2} & 0 & 0 & 0 & 0 & 3c_1 \\
 m_{2 \rightarrow 1} & \cdot & 0 & 0 & m_{1 \rightarrow 2} & m_{1 \rightarrow 2} & 0 & 0 & c_1 \\
 m_{2 \rightarrow 1} & 0 & \cdot & 0 & m_{1 \rightarrow 2} & 0 & m_{1 \rightarrow 2} & 0 & c_1 \\
 m_{2 \rightarrow 1} & 0 & 0 & \cdot & 0 & m_{1 \rightarrow 2} & m_{1 \rightarrow 2} & 0 & c_1 \\
 0 & m_{2 \rightarrow 1} & m_{2 \rightarrow 1} & 0 & \cdot & 0 & 0 & m_{1 \rightarrow 2} & c_2 \\
 0 & m_{2 \rightarrow 1} & 0 & m_{2 \rightarrow 1} & 0 & \cdot & 0 & m_{1 \rightarrow 2} & c_2 \\
 0 & 0 & m_{2 \rightarrow 1} & m_{2 \rightarrow 1} & 0 & 0 & \cdot & m_{1 \rightarrow 2} & c_2 \\
 0 & 0 & 0 & 0 & m_{2 \rightarrow 1} & m_{2 \rightarrow 1} & m_{2 \rightarrow 1} & \cdot & 3c_2 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0
 \end{pmatrix}$$

We wish to calculate:

$$\mathbb{P}(X_t = (1 \mid 2\ 3) \mid X_0 = (1\ 2\ 3 \mid)) = \mathbb{P}_\pi(X_t = 1 \mid 2\ 3) = \left(e^{Qt} \right)_{1\ 2\ 3 \mid, 1 \mid 2\ 3}$$

$$\tilde{f} : 1 \mapsto 0, 2 \mapsto 0, 3 \mapsto 0$$

Applying lumpability

\hat{Q} :

$$\begin{array}{l}
 000 | \\
 00 | 0 \\
 0 | 00 \\
 | 000 \\
 A
 \end{array}
 \left(
 \begin{array}{c|ccc|c}
 000 | & 00 | 0 & 0 | 00 & | 000 & A \\
 \cdot & 3m_{1 \rightarrow 2} & 0 & 0 & 3c_1 \\
 m_{2 \rightarrow 1} & \cdot & 2m_{1 \rightarrow 2} & 0 & c_1 \\
 0 & 2m_{2 \rightarrow 1} & \cdot & 0 & c_2 \\
 0 & 0 & 2m_{2 \rightarrow 1} & \cdot & 3c_2 \\
 0 & 0 & 0 & 0 & 0
 \end{array}
 \right)$$

$$\mathbb{P}_\pi(X_t = (1 | 23)) = \frac{1}{3} \mathbb{P}_{\hat{\pi}}(\hat{X}_t = (0 | 00)) = \frac{1}{3} \left(e^{\hat{Q}t} \right)_{000 |, 0 | 00}$$

Reduced statespace

L	2	3	4	5	6
States	6	22	94	454	2430

Table : The number of states for $L = 2, 3, 4, 5, 6$ lineages

L	2	3	4	5	6
States	7	17	39	67	117

Table : The number of states in the reduced states-space for $L = 2, 3, 4, 5, 6$ lineages.

Thank you.

