

5. Modèles de dépendences

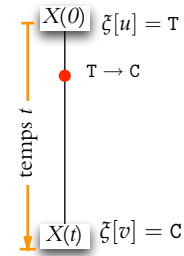
Procéssus de Markov à temps continu

Le PMTC est un procéssus $X(t)$ qui décrit les variables aléatoires au temps $t \in [0, \infty)$. Dans le contexte d'ADN, l'évolution entre le nœud parent u et l'enfant v suit un PMTC $X(t) \in \{A, C, G, T\}$: on lance le procéssus avec $X(0) = \zeta[u]$ (l'étiquette de u), et on met $\zeta[v] = X(t)$ où t dénote la longueur de l'arête uv . On a alors la matrice de probabilités conditionnelles

$$\begin{aligned} \mathbf{M}(t) &= \left[\mathbb{P} \left\{ \zeta[v] = j \mid \zeta[u] = i \right\} : i, j = A, C, G, T \right] && (\text{substitutions } i \rightarrow j \text{ sur arête } uv) \\ &= \begin{bmatrix} p_{A \rightarrow A}(t) & p_{A \rightarrow C}(t) & p_{A \rightarrow G}(t) & p_{A \rightarrow T}(t) \\ p_{C \rightarrow A}(t) & p_{C \rightarrow C}(t) & p_{C \rightarrow G}(t) & p_{C \rightarrow T}(t) \\ p_{G \rightarrow A}(t) & p_{G \rightarrow C}(t) & p_{G \rightarrow G}(t) & p_{G \rightarrow T}(t) \\ p_{T \rightarrow A}(t) & p_{T \rightarrow C}(t) & p_{T \rightarrow G}(t) & p_{T \rightarrow T}(t) \end{bmatrix} \\ &= \exp(\mathbf{Q}t) && (\mathbf{Q} \text{ est la matrice de taux instantané}) \end{aligned}$$

où l'exponentiation¹ matricielle dénote la somme infinie

$$\exp(\mathbf{Q}t) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} + \mathbf{Q}t + \frac{(\mathbf{Q}t)^2}{2} + \frac{(\mathbf{Q}t)^3}{3!} + \dots = \sum_{k=0}^{\infty} \frac{(\mathbf{Q}t)^k}{k!}.$$



¹ à l'analogie de la série Taylor $e^{qt} = \sum_{k=0}^{\infty} \frac{(qt)^k}{k!}$

Modèle de Jukes-Cantor. La matrice de taux

$$\mathbf{Q} = \begin{bmatrix} -1 & 1/3 & 1/3 & 1/3 \\ 1/3 & -1 & 1/3 & 1/3 \\ 1/3 & 1/3 & -1 & 1/3 \\ 1/3 & 1/3 & 1/3 & -1 \end{bmatrix} \quad \text{donne} \quad p_{i \rightarrow j}(t) = \begin{cases} \frac{1}{4} + \frac{3}{4}e^{-4t/3} & \{i = j\} \\ \frac{1}{4} - \frac{1}{4}e^{-4t/3} & \{i \neq j\} \end{cases}$$

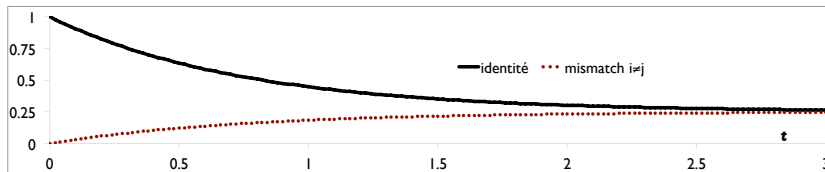


FIG. 1: Probabilités $p_{i \rightarrow j}(t)$ dans le modèle Jukes-Cantor.

Noter qu'à $t = 0$ (arête de longueur nulle), on a

$$\mathbb{P} \left\{ X(0) = i \mid X(0) = i \right\} = 1, \text{ et}$$

$$\lim_{t \rightarrow \infty} p_{i \rightarrow j}(t) = \lim_{t \rightarrow \infty} \mathbb{P} \left\{ X(t) = j \mid X(0) = i \right\} = 1/4 \quad \text{pour tout } i, j = A, C, G, T.$$

Modèle F84/HKY. Le modèle de Felsenstein (1984) ou Hasegawa-Kishino-Yano (1985) s'écrit avec une matrice de taux

$$Q = \begin{bmatrix} \cdot & \pi_C \alpha & \pi_G \beta & \pi_T \alpha \\ \pi_A \alpha & \cdot & \pi_G \alpha & \pi_T \beta \\ \pi_A \beta & \pi_C \alpha & \cdot & \pi_T \alpha \\ \pi_A \alpha & \pi_C \beta & \pi_G \alpha & \cdot \end{bmatrix} \quad (\text{somme est } 0 \text{ dans toute rangée})$$

défini mène à la composition stationnaire $\lim_{t \rightarrow \infty} \mathbb{P}\{X(t) = i\} = \pi_i$ et un rapport de taux transition : transversion $\kappa = \frac{\beta}{\alpha}$.

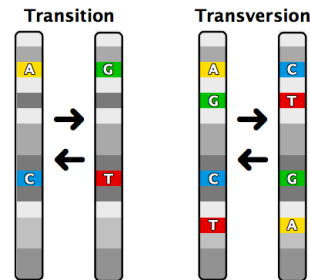
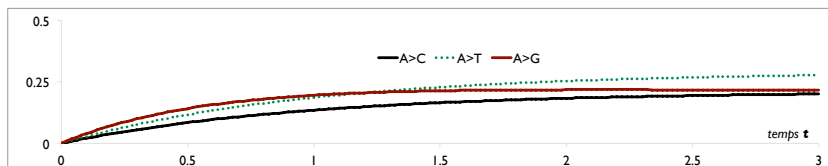


FIG. 2: Transitions et transversions
 FIG. 3: Probabilités $p_{i \rightarrow j}(t)$ dans le modèle F84 pour GC% = 42% ($\pi_A = p_{i_T} = 0.29$, $\pi_C = \pi_G = 0.21$) et $\kappa = 2$. On a $p_{A \rightarrow C} = 0.21(1 - e^{-t/0.9744})$, $p_{A \rightarrow T} = 0.29(1 - e^{-t/0.9744})$, et $p_{A \rightarrow G} = 0.21(1 + e^{-t/0.9744} - 2e^{-t/0.6496})$.

Multiplication. On a que $\mathbf{M}(t) \cdot \mathbf{M}(s) = \mathbf{M}(t+s)$ pour tout $s, t \geq 0$, et que $\boldsymbol{\pi} \cdot \mathbf{M}(t) = \boldsymbol{\pi}$ avec la distribution stationnaire $\boldsymbol{\pi}$.

Segmentation

Segmentation par blocs maximaux. On a une séquence de scores $a[1..\ell]$ avec $a[j] \in \mathbb{R}$ (peut être négatif ou positif). On pénalise un bloc $j..j+k-1$ par la somme $a[j] + a[j+1] + \dots + a[j+k-1]$. Bloc de score maximal par programmation dynamique : $S \leftarrow 0$, et faire $S \leftarrow \max\{0, S + a[j]\}$ avec $j = 1, 2, \dots, \ell$.

Modèle de Markov caché. Chaque colonne $k = 1, \dots, \ell$ correspond à un état caché $q[k] \in \mathcal{Q}$: par exemple, $\mathcal{Q} = \{\text{neutre}, \text{conservé}\}$. Émission = résidues alignées $x[k] = (x_1[k] \dots x_n[k])$ dans la colonne. Le modèle phylogénétique définit la probabilité $p(i, x[k])$ pour la colonne est émise dans l'état i . La séquence d'états $q[1..\ell]$ est au hasard avec $\tau(i, j) = \mathbb{P}\{q[k+1] = j \mid q[j] = i\}$. Quel est la séquence d'états la plus probable ou le **chemin Viterbi**? Pour cela on définit la vraisemblance du meilleur préfixe, ce qu'on peut calculer par programmation dynamique :

$$\delta_1(i) = \pi(i) \cdot p(i, x[1]) \quad (5.1a)$$

$$\delta_k(i) = \left(\max_{j \in \mathcal{Q}} \delta_{k-1}(j) \cdot \tau(j, i) \right) \cdot p(i, x[k]) \quad \{k > 0\} \quad (5.1b)$$

On reconstruit le *chemin Viterbi* qui correspond à $\max_i \delta_\ell(i)$ par *backtracking* : il faut stocker le meilleur état qui précède j , ainsi que $\delta_k(i)$ dans la récurrence de (5.1b).

p.e., $a[k] = \text{LODS}$ pour deux hypothèses.