

3. Substitutions

Modèle iid. Modèle probabiliste pour une séquence : caractères aléatoires iid¹ : $\mathbb{P}\{S[i] = a\} = \pi_a$. Par conséquent, le nombre d'occurrences $\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}$ détermine la probabilité d'une séquence, sans égard à l'ordre des nucléotides²

Pondération de substitutions

Supposons qu'on a deux séquences ADN $s[0..\ell - 1]$ et $t[0..\ell - 1]$, qui correspondent à une région sans trous dans un alignement. On veut décider s'il y a une dépendance entre les résidus alignés. On compare deux modèles paramétriques, correspondant à hypothèses distincts. Les modèles déterminent la distribution de deux séquences aléatoires $S[0..\ell - 1]$ et $T[0..\ell - 1]$, comprenant des paires iid $(S[i], T[i])$. Ici, S et T sont des variables aléatoires. On observe s et t : les modèles définissent $\mathbb{P}\{S = s; T = t\}$.

Hypothèse 0 : S et T sont indépendantes. On définit

$$L_0 = \mathbb{P}\{S = s; T = t \mid H_0\} = \prod_{i=0}^{\ell-1} \pi_{s[i]} \cdot \pi_{t[i]}. \quad (3.1a)$$

Hypothèse 1 : S et T sont dépendantes :

$$L_1 = \mathbb{P}\{S = s; T = t \mid H_1\} = \prod_{i=0}^{\ell-1} \pi_{s[i]} \cdot \mathbb{P}\{T[i] = t[i] \mid S[i] = s[i]\}. \quad (3.1b)$$

Les quantités L_0, L_1 mesurent le support aux données par les hypothèses. En général, $\mathbb{P}\{\text{données} \mid \text{hypothèse}\}$ s'appelle la **vraisemblance**³

Matrice de substitution. La **matrice de substitution** \mathbf{M} se définit par la distribution conditionnelle

$$\mathbf{M}_{a,b} = \mathbb{P}\{T[i] = b \mid S[i] = a\}$$

(à n'importe quelle position i selon la supposition iid). C'est une **matrice stochastique** de taille 4×4 : $0 \leq \mathbf{M}_{a,b} \leq 1$ pour tout $a, b = \mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}$, ainsi que $\sum_b \mathbf{M}_{a,b} = 1$ pour tout a .

Rapport des chances. Le rapport des vraisemblances de (3.1) s'appelle le **rapport des chances** (*odds ratio*). Le *log-odds-ratio* quantifie l'évidence pour hypothèse 1 par rapport à hypothèse 0 sur une échelle logarithmique⁴

$$\text{LODS} = \log \frac{L_1}{L_0}. \quad (3.2)$$

Par (3.1a) et (3.1b), on a

$$\text{LODS} = \log \frac{L_1}{L_0} = \log \frac{\prod_{i=0}^{\ell-1} \pi_{s[i]} \cdot \mathbf{M}_{s[i],t[i]}}{\prod_{i=0}^{\ell-1} \pi_{s[i]} \cdot \pi_{t[i]}} = \sum_{i=0}^{\ell-1} \log \frac{\mathbf{M}_{s[i],t[i]}}{\pi_{t[i]}}.$$

¹ iid=indépendants et identiquement distribués

² Exemple : $\pi_{\mathbf{A}} = \pi_{\mathbf{T}} = 32\%$, $\pi_{\mathbf{C}} = \pi_{\mathbf{G}} = 18\%$, ou taux de $(\mathbf{G} + \mathbf{C})$ à 36%. On a alors $\mathbb{P}\{S = \mathbf{ACAT}\} = \mathbb{P}\{S = \mathbf{CTAA}\} = 0.32^3 \times 0.18 = 0.00589 \dots$

³ Noter que c'est tout à fait différent de $\mathbb{P}\{\text{hypothèse} \mid \text{données}\}$.

⁴ On mesure l'évidence en *bits* par $\lg = \log_2$, en *nats* par $\ln = \log_e$, et en *ban* par \log_{10} . L'échelle $10 \log_{10} = \log_{1.2589\dots}$ («Phred scale» ou deciban) est utilisé en séquençage.

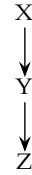
Matrice de pondération. On définit la pondération de paires (a, b) alignées par la valeur arrondie

$$C_{a,b} = \left\lceil \log_{\alpha} \frac{\mathbf{M}_{a,b}}{\pi_b} \right\rceil,$$

sur une échelle logarithmique quelconque $\alpha > 1$. L'évidence pour hypothèse 1 est la somme de poids : $\text{LODS} = \sum_{i=0}^{\ell-1} C_{s[i],t[i]}$.

Substitutions multiples. Évolution de caractères homologues est **sans mémoire** : états successifs forment une chaîne de Markov. Probabilités de **substitution** selon la matrice $\mathbf{M}_{X \rightarrow Y}$ où $\mathbf{M}_{X \rightarrow Y}[a, b] = \mathbb{P}\{Y = b \mid X = a\}$ avec valeurs $a, b \in \Sigma = \{1, \dots, r\}$. (Exemple : $\Sigma = \{1, 2, 3, 4\}$ encodant ADN)

$$\begin{aligned} \mathbb{P}\{Z = z \mid X = x\} &= \sum_{y=1}^r \mathbb{P}\{Z = z, Y = y \mid X = x\} \\ &= \sum_{y=1}^r \mathbb{P}\{Z = z \mid X = x, Y = y\} \mathbb{P}\{Y = y \mid X = x\} \\ &= \sum_{y=1}^r \mathbb{P}\{Z = z \mid Y = y\} \mathbb{P}\{Y = y \mid X = x\} \quad (\text{propriété de Markov}) \end{aligned}$$



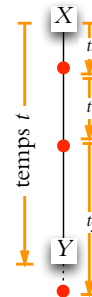
Donc, on peut écrire que $\mathbf{M}_{X \rightarrow Z} = \mathbf{M}_{X \rightarrow Y} \cdot \mathbf{M}_{Y \rightarrow Z}$

Temps continu. À chaque mutation, le changement d'état est déterminé par la matrice \mathbf{M} . Si k événements en temps t , alors $\mathbf{M}_{X \rightarrow Y} = \mathbf{M}^k$. Donc la matrice pour temps t devrait être $\mathbf{M}^{N(t)}$ où $N(t)$ est le nombre de mutations pendant temps t . Problème : $N(t)$ est aléatoire. . .

☞ **«orderly»** La probabilité qu'il y ait plus d'une occurrence dans un petit intervalle de temps est négligeable

$$\lim_{\delta \rightarrow 0} \mathbb{P}\{N(t + \delta) > 1 \mid N(t + \delta) \geq 1\} = 0$$

☞ **«memoryless»** Le nombre d'occurrences pendant un intervalle ne dépend pas des événements précédents : $N(t + s) - N(t)$ est indépendant de $N(t)$



Processus de Poisson. Si le processus est *orderly* et *memoryless*, alors c'est un **processus de Poisson** avec les propriétés suivantes.

- ★ nombre d'arrivées pendant temps t est une v.a. Poisson : $\mathbb{P}\{N(t + s) - N(s) = k\} = e^{-\lambda t} \frac{(\lambda t)^k}{k!}$. Le paramètre λ est l'intensité ou le **taux** du processus.
- ★ temps d'attente est une v.a. exponentielle : $\forall t, s : \mathbb{P}\{N(s + t) = N(s)\} = e^{-\lambda t}$

Processus de Markov. En retournant à ce qui se passe pendant temps t :

$$\sum_{k=0}^{\infty} e^{-\lambda t} \frac{(\lambda t)^k}{k!} \mathbf{M}^k = \exp(\lambda t(\mathbf{M} - \mathbf{I})) = \exp(\lambda t \mathbf{Q})$$

où \mathbf{I} est la matrice d'unité. $\mathbf{Q} = \mathbf{M} - \mathbf{I}$ est la *matrice instantanée de taux de substitutions*. Calcul de $\exp(\lambda \mathbf{Q} t)$: décomposition de la matrice $\mathbf{Q} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}$, et faire $e^{\lambda \mathbf{Q} t} = \mathbf{U} e^{\lambda t \mathbf{\Lambda}} \mathbf{V}$.