

4. Variables inobservées

Génomique comparative. Principe de génomique comparative : sélection négative ralentit le taux d'évolution \Rightarrow plus de conservation de séquence à un élément fonctionnel.

Inférence par logarithme des chances à partir d'un alignement : hypothèse 0 = substitutions par $e^{\mu Q^t}$, hypothèse 1 = substitutions par $e^{\rho \mu Q^t}$, on connaît μt , Q , $\rho < 1$. Alors, on calcule les vraisemblances L_0, L_1 selon les hypothèses et utilise $LLR = \log \frac{L_1}{L_0}$ pour détecter l'**ombre** de sélection.

Avec plus de séquences, on peut détecter la conservation mieux. . . mais comment calculer L_i ?

Vraisemblance sur une phylogénie

Modèle markovien pour étiquetage (v.a. $\zeta[u] \in \mathcal{X}$ à tout nœud u) :

- (1) distribution à la racine par $\pi_x = \mathbb{P}\{\zeta[\text{racine}] = x\}$;
- (2) probabilités de changement d'étiquette sur arête uv

$$p_{x \rightarrow y}(uv) = \mathbb{P}\{\zeta[v] = y \mid \zeta[u] = x\}$$

Vraisemblance pour un étiquetage x des nœuds terminaux — il faut considérer tous les étiquetages de tous les nœuds

$$L[x] = \sum_{\tilde{x}: \text{étiquettes possibles aux nœuds}} \pi_{\tilde{x}[\text{racine}]} \prod_{\text{arêtes } uv} p_{\tilde{x}[u] \rightarrow \tilde{x}[v]}(uv).$$

Calcul efficace par programmation dynamique (algorithme de Felsenstein — *peeling*)

Idée : définir la vraisemblance conditionnelle $L_u[x]$ pour tout caractère $x \in \mathcal{X}$: probabilité pour étiquettes dans le sous-arbre de u quand $\zeta[u] = x$.

$$L_u[x] = \{\zeta[u] = x\} \quad u \text{ terminal}$$

$$L_u[x] = \prod_{v \in \text{enfants}(u)} \left(\sum_{y \in \mathcal{X}} p_{x \rightarrow y}(uv) L_v[y] \right) \quad u \text{ ancestral}$$

Vraisemblance et modèles graphiques

On peut généraliser l'approche à des graphes orientés. Le modèle probabiliste comprend un ensemble de variables aléatoires inobservées Y_1, \dots, Y_m et l'observation Y_0 même :

$$p(x|\theta) = \sum_{y_1, y_2, \dots, y_m} \mathbb{P}\{Y_0 = x \mid Y_1 = y_1, \dots, Y_m = y_m; \theta\} \cdot \mathbb{P}\{Y_1 = y_1, \dots, Y_m = y_m \mid \theta\}.$$

Ici, θ représente les paramètres du modèle. La somme suggère une temps exponentiel en m , mais on peut souvent exploiter les indépendances conditionnelles entre les Y_i pour accélérer le calcul. Considérons le graphe G de

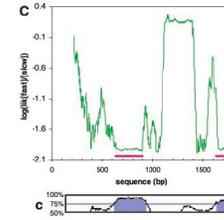
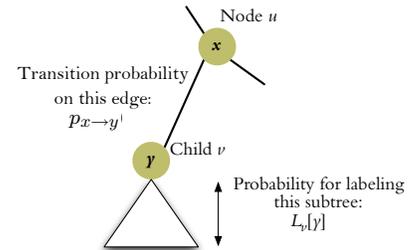


FIG. 1: **Phylogenetic shadowing** illustré par Boffelli & al. *Science* 299 :1391 (2003). En haut : LLR dans un fenêtre de 50 pb glissant dans l'alignement multiple de 5 primates, avec des exons indiqués. En bas : % identité humain-souris.



dépendances entre les variables aléatoires. Les nœuds de G correspondent aux variables $Y_i: i = 0, 1, \dots, m$. Les arcs de G sont orientés. L'ensemble des arcs entrants au nœud j , dénoté par $\text{pred}(j) = \{i: ij \in G\}$, définit l'indépendance conditionnelle :

$$\mathbb{P}\{Y_j \mid Y_i: i \neq j\} = \mathbb{P}\{Y_j \mid Y_i: i \in \text{pred}(j)\}.$$

Un *modèle graphique* est un modèle probabiliste dont le graphe de dépendances ne contient aucun cycle.

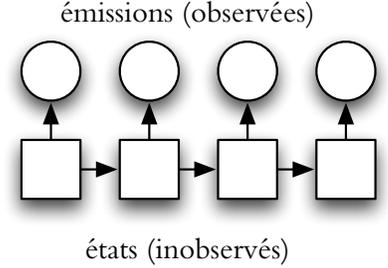
Modèle de Markov caché

Un **modèle de Markov caché** (*hidden Markov model* — HMM) sur un alphabet \mathcal{A} comprend les paramètres suivants.

- ★ ensemble d'états \mathcal{Q} de taille $|\mathcal{Q}| = N$
- ★ probabilités de transition entre états : $\tau: \mathcal{Q} \times \mathcal{Q} \rightarrow [0, 1]$
- ★ probabilités d'émission à chaque état : $p: \mathcal{Q} \times \mathcal{A} \rightarrow [0, 1]$
- ★ probabilités d'état initial : $\pi: \mathcal{Q} \rightarrow [0, 1]$

Le modèle $M = (\mathcal{Q}, \tau, p, \pi)$ définit une distribution sur séquences d'une longueur fixe. En particulier, M produit une séquence $X[0..\ell - 1]$ selon la procédure aléatoire suivante.

G1 choisir l'état initial $q[0]$ au hasard selon les probabilités π
 G2 **for** $i \leftarrow 0, 1, \dots, \ell - 1$
 G3 choisir $X[i]$ au hasard par les probabilités d'émission $p(q[i], \cdot)$
 G4 choisir prochain état $q[i + 1]$ au hasard par les probabilités de transition $p(q[i], \cdot)$



Vraisemblance. La vraisemblance selon M se définit par

$$L(M) = \mathbb{P}\{X = x \mid M\} = \sum_{q[0..\ell-1] \in \mathcal{Q}^\ell} \left(\pi(q[0]) \cdot p(q[0], x[0]) \right. \quad (1er \text{ état+symbole}) \\ \times \tau(q[0], q[1]) \cdot p(q[1], x[1]) \times \dots \quad (2e \text{ état+symbole}) \\ \left. \times \tau(q[\ell - 2], q[\ell - 1]) \cdot p(q[\ell - 1], x[\ell - 1]) \right). \quad (4.1)$$

Probabilité de préfixe. On définit la vraisemblance sur les préfixes (*forward probabilities*)

$$\alpha_k(i) = \sum_{q[0..k]; q[k]=i} \mathbb{P}\{X[0..k] = x[0..k] \mid M\},$$

ce qu'on peut calculer par programmation dynamique.

$$\alpha_0(q) = \pi(q) \quad (4.2a)$$

$$\alpha_k(i) = \left(\sum_{j \in \mathcal{Q}} \alpha_{k-1}(j) \cdot \tau(j, i) \right) \cdot p(i, x[k]) \quad \{k > 0\} \quad (4.2b)$$