

# Un modèle d'évolution des séquences d'ADN avec dépendance au voisin de gauche

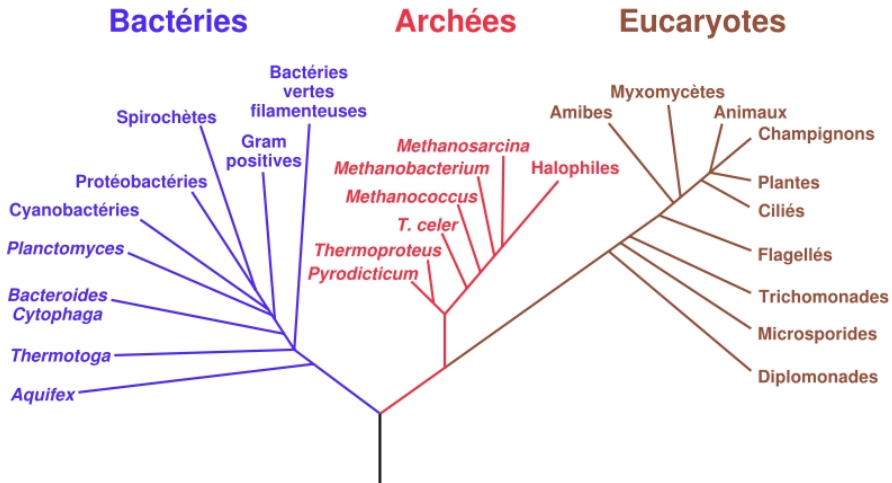
Audrey Finkler

Colloque Jeunes Probabilistes et Statisticiens, Mont-Dore

3 mai 2010



# L'arbre phylogénétique des vivants



Distance évolutive :  $K = kT$ .

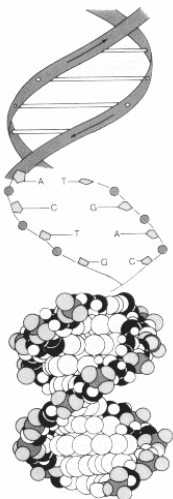
# Plan

- 1 Généralités sur les modèles d'évolution
- 2 Modèle avec dépendance au voisin de gauche
- 3 Simulations et applications

# Plan

- 1 Généralités sur les modèles d'évolution
- 2 Modèle avec dépendance au voisin de gauche
- 3 Simulations et applications

## Quelques notions de biologie moléculaire



- Séquence d'ADN :

$$\mathbf{X} = (X_1, \dots, X_n),$$

suite finie de variables aléatoires avec  
 $X_i \in \mathcal{A} = \{A, C, G, T\}$ .

- Complémentarité des bases :

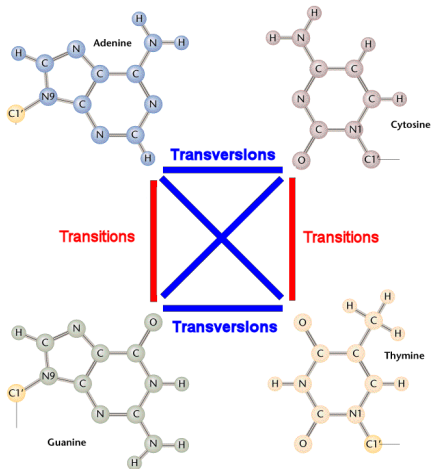


- Purines :  $\mathcal{R} = \{A, G\}$ .  
Pyrimidines :  $\mathcal{Y} = \{C, T\}$ .

# Processus de substitution markovien à temps continu :

- substitution : mutation ponctuelle sans insertion ni délétion,
- temps continu,
- markovien d'ordre 1,
- homogène, uniforme et irréductible,
- sites indépendants.

# Transitions et transversions



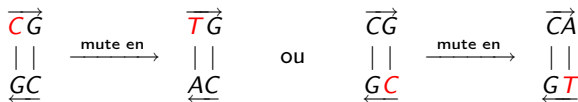
- Histoire du site  $i$  :

$$X_i := (X_i(t))_{t \in [0, \tau]}.$$

- **Transition** (ts) :  $\mathcal{R} \rightarrow \mathcal{R}$  ou  $\mathcal{Y} \rightarrow \mathcal{Y}$ .
- **Transversion** (tv) :  $\mathcal{R} \rightarrow \mathcal{Y}$  ou  $\mathcal{Y} \rightarrow \mathcal{R}$ .

- $1 < \frac{ts}{tv} < 10$ .

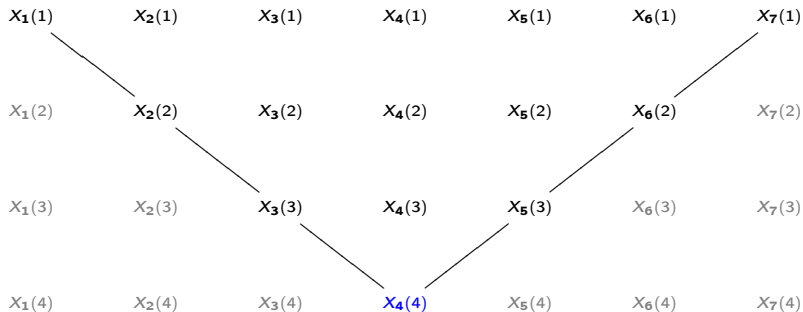
## Effet CpG



Substitutions à partir d'un dinucléotide CpG sur un brin et son complémentaire.



# Cône de dépendances

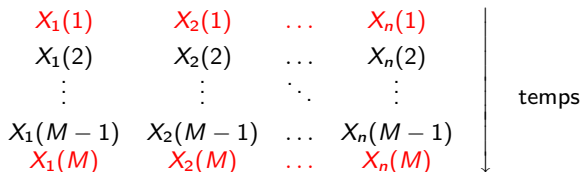


# Plan

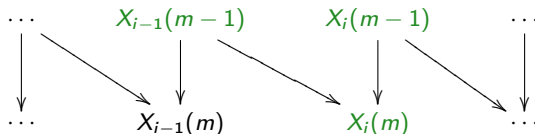
- 1 Généralités sur les modèles d'évolution
- 2 **Modèle avec dépendance au voisin de gauche**
- 3 Simulations et applications

# Modèle

Temps discrétisé :

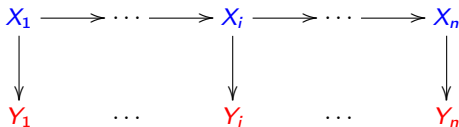


Graphe des dépendances du modèle :



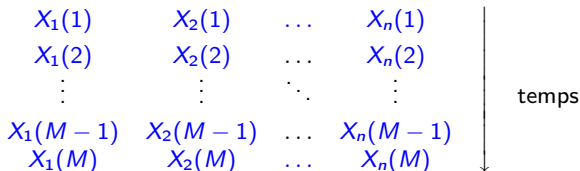
# Modèle vu comme une chaîne de Markov cachée (HMM)

Deux processus emboîtés :



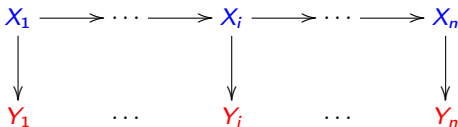
Pour  $i \in \{1, \dots, n\}$ ,

- états **cachés**  $X_i = (X_i(m))_{1 \leq m \leq M}$ ,
- états **observés**  $Y_i = \{X_i(1), X_i(M)\}$ .



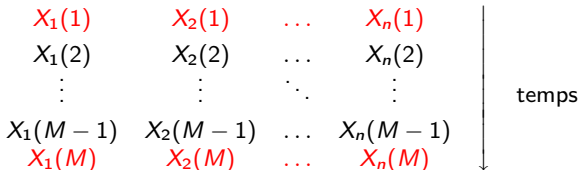
# Modèle vu comme une chaîne de Markov cachée (HMM)

Deux processus emboîtés :



Pour  $i \in \{1, \dots, n\}$ ,

- états cachés  $X_i = (X_i(m))_{1 \leq m \leq M}$ ,
- états **observés**  $Y_i = \{X_i(1), X_i(M)\}$ .





## Paramètres à estimer

Posons :

$$c_{\theta}(u; v, w) = \begin{cases} \mu \cdot \tau^{\epsilon(u,v)} \cdot \left(\frac{1-\tau}{2}\right)^{1-\epsilon(u,v)} \cdot \gamma^{\kappa(w,v)} & \text{si } v \neq u, \\ 1 - \mu \cdot \gamma^{\kappa(w,v)} & \text{si } v = u, \end{cases}$$

où  $\theta = (\mu, \tau, \gamma) \in [0, 1]^3$  avec

- $\mu$  : probabilité pour un site de subir une substitution,
- $\tau$  : facteur multiplicatif  $T_s/T_v$ , avec  $\epsilon(u, v)$  indicatrice d'une transition entre  $u$  et  $v$ ,
- $\gamma$  : facteur multiplicatif prenant en compte le contexte,  $\kappa(w, v)$  indiquant un contexte différent de  $CpG$  ou  $TpA$ .

## Estimation de $\theta$ par l'algorithme EM-HMM

$$Q(\theta, \theta_n^{(k)}) = q + \sum_{i=2}^n \sum_{U \in \mathcal{A}^M} \sum_{V \in \mathcal{A}^M} \mathbb{P}_k(X_{i-1} = U, X_i = V | Y) \ln D(U, V).$$

On approche itérativement l'estimateur du maximum de vraisemblance de  $\theta$  par l'algorithme EM appliqué aux HMM. Le paramètre est initialisé à  $\theta_n^{(0)}$  puis on alterne les étapes suivantes :

**E** : calculer  $Q(\theta, \theta_n^{(k)})$ ,

**M** : poser  $\theta_n^{(k+1)} = \arg \max_{\theta} Q(\theta, \theta_n^{(k)})$ .



# Plan

- 1 Généralités sur les modèles d'évolution
- 2 Modèle avec dépendance au voisin de gauche
- 3 Simulations et applications

# Simulations

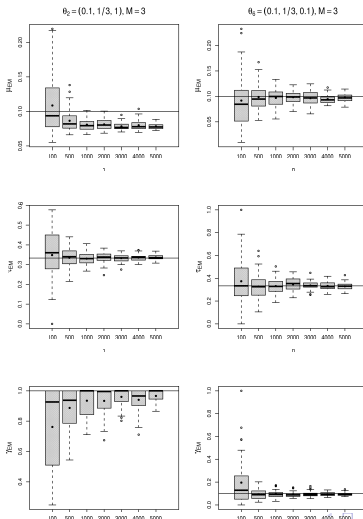
Pour  $n \in \{100, 500, 1\ 000, 2\ 000, 3\ 000, 4\ 000, 5\ 000\}$  :

- 50 répétitions pour  $M = 3, 4$ ,
- 25 répétitions pour  $M = 5$ ,

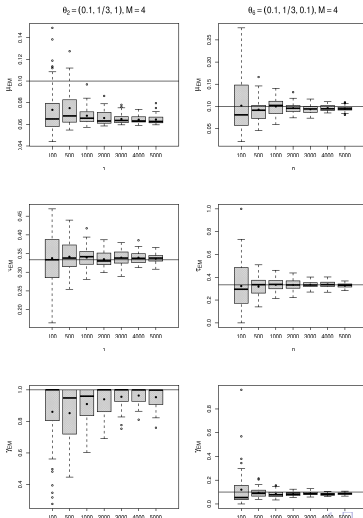
pour chacune des sept valeurs de  $\theta$  :

	$\mu$	$\tau$	$\gamma$
$\theta_1$	0.01	1/3	1
$\theta_2$	0.1	1/3	1
$\theta_3$	0.1	0.5	1
$\theta_4$	0.1	0.8	1
$\theta_5$	0.1	1/3	1/3
$\theta_6$	0.1	1/3	0.1
$\theta_7$	0.01	0.8	1/3

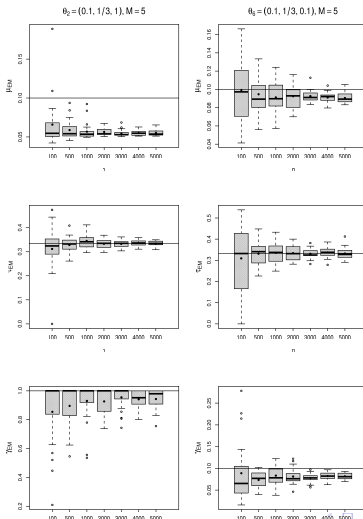
# Résultats pour $M = 3$



# Résultats pour $M = 4$



# Résultats pour $M = 5$



## Application à des données réelles

Le taux global de substitution  $k_{EM}$  est tel que  $\mu_{EM} = \frac{k_{EM} T}{M - 1}$ .

- séquence codante de l'insuline : **Souris** et **Rat** pour  $n = 327$  et  $M = 5$  :

$$\mu_{EM} = 0.060, \quad \tau_{EM} = 0.717, \quad \gamma_{EM} = 0.241.$$

Durée depuis divergence : 41 Ma, soit  $k_{EM} = 0.3\%$  par site et Ma.

- séquence non-codante de la bêta-globine : **Homme** et **Chimpanzé** pour  $n = 2\,147$  et  $M = 5$  :

$$\mu_{EM} = 0.004, \quad \tau_{EM} = 0.634, \quad \gamma_{EM} = 0.839.$$









Durée depuis divergence : 4.7 Ma, soit  $k_{EM} = 0.17\%$  par site et Ma.

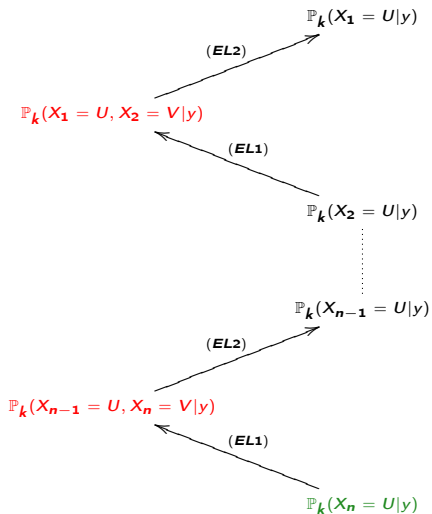
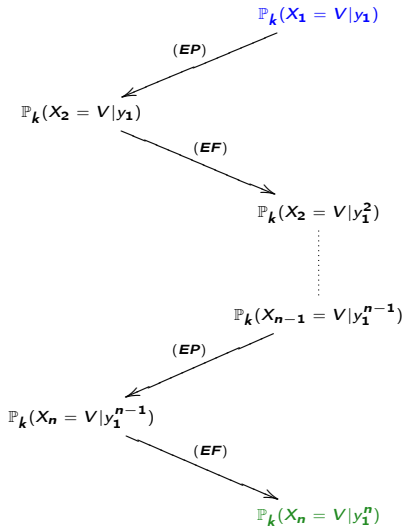
## Généralisations possibles

- Considérer plus de deux séquences : famille de séquences alignées.
- Remplacer les nucléotides par des codons.
- Étendre la dépendance au voisin de droite : champs markoviens.

Merci pour votre attention !



-  L. Baum, T. Petrie, Statistical inference for probabilistic functions of finite state Markov chains [Ann. Math. Stat.](#), 37 :1554-1563, 1966
-  L. Baum, T. Petrie, G. Soules, N. Weiss, A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains [Ann. Math. Stat.](#), 41(1) :164-171, 1970
-  M. Bulmer, Neighboring base effects on substitution rates in pseudogenes [Mol. Biol. and Evol.](#), 3 :322-329, 1986
-  J.L. Jensen, Context dependent DNA evolutionary models [Research Reports No.458](#), May 2005
-  J.L. Jensen, A.-M. Pedersen, Probabilistic models of DNA sequence evolution with context dependent rates of substitution [Adv. Appl. Probab.](#), 32 :499-517, 2000
-  A.-M. Pedersen, J.L. Jensen, A dependent-rates model and an MCMC-based methodology for the maximum-likelihood analysis of sequences with overlapping reading frames [Mol. Biol. Evol.](#), 18(5) :763-776, 2001
-  A.-M. Pedersen, C. Wiuf, F.B. Christiansen, A codon-based model designed to describe lentiviral evolution [Mol. Biol. Evol.](#), 15(8) :1069-1081, 1998
-  C. Wu, On the convergence properties of the EM algorithm [Ann. Stat.](#), 11(1) :95-103., 1983



# Equations prédictives (EP)

Notons  $Y_i^j = (Y_i, \dots, Y_j)$ .

$$\begin{aligned}
 \mathbb{P}_k(X_i = V | Y_1^{i-1}) &= \sum_{U \in \mathcal{A}^M} \mathbb{P}_k(X_{i-1} = U, X_i = V | Y_1^{i-1}), \\
 &= \sum_{U \in \mathcal{A}^M} \mathbb{P}_k(X_i = V | X_{i-1} = U, Y_1^{i-1}) \mathbb{P}_k(X_{i-1} = U | Y_1^{i-1}), \\
 &= \sum_{U \in \mathcal{A}^M} \mathbb{P}_k(X_i = V | X_{i-1} = U) \mathbb{1}_{(u_1, u_M) = y_{i-1}} \mathbb{P}_k(X_{i-1} = U | Y_1^{i-1}), \\
 &= \sum_{U \in \mathcal{A}^M} C^{(k)}(U, V) \mathbb{P}_k(X_{i-1} = U | Y_1^{i-1}) \mathbb{1}_{(u_1, u_M) = y_{i-1}}.
 \end{aligned}$$

## Equations de filtrage (EF)

$$\begin{aligned}\mathbb{P}_k(X_i = U | y_1^i) &= \mathbb{P}_k(X_i = U | y_1^{i-1}, y_i), \\ &= \frac{\mathbb{P}_k(X_i = U, y_i | y_1^{i-1})}{\mathbb{P}_k(y_i | y_1^{i-1})}, \\ &= \frac{\mathbb{P}_k(X_i = U, y_i | y_1^{i-1})}{\sum_{V \in \mathcal{A}^M} \mathbb{P}_k(X_i = V, y_i | y_1^{i-1})}, \\ &= \frac{\mathbb{P}_k(y_i | X_i = U, y_1^{i-1}) \mathbb{P}_k(X_i = U | y_1^{i-1})}{\sum_{V \in \mathcal{A}^M} \mathbb{P}_k(y_i | X_i = V, y_1^{i-1}) \mathbb{P}_k(X_i = V | y_1^{i-1})}, \\ &= \frac{\mathbb{1}_{(u_1, u_M) = y_i} \mathbb{P}_k(X_i = U | y_1^{i-1})}{\sum_{V \in \mathcal{A}^M} \mathbb{1}_{(v_1, v_M) = y_i} \mathbb{P}_k(X_i = V | y_1^{i-1})}.\end{aligned}$$

# Equations de lissage 1 (EL1)

$$\begin{aligned}
 \mathbb{P}_k(X_{i-1} = U, X_i = V|y) &= \mathbb{P}_k(X_{i-1} = U|X_i = V, y)\mathbb{P}_k(X_i = V|y), \\
 &= \mathbb{P}_k(X_{i-1} = U|X_i = V, y_1^{i-1})\mathbb{P}_k(X_i = V|y)\mathbb{1}_{(\mathbf{v}_0, \mathbf{v}_M)=y_i}, \\
 &= \frac{\mathbb{P}_k(X_i = V, X_{i-1} = U|y_1^{i-1})\mathbb{P}_k(X_i = V|y)}{\mathbb{P}_k(X_i = V|y_1^{i-1})}\mathbb{1}_{(\mathbf{v}_1, \mathbf{v}_M)=y_i}, \\
 &= \frac{\mathbb{P}_k(X_i = V|X_{i-1} = U, y_1^{i-1})\mathbb{P}_k(X_{i-1} = U|y_1^{i-1})}{\mathbb{P}_k(X_i = V|y_1^{i-1})} \cdot \\
 &\quad \mathbb{P}_k(X_i = V|y)\mathbb{1}_{(\mathbf{v}_1, \mathbf{v}_M)=y_i}, \\
 &= \frac{\mathbb{P}_k(X_i = V|X_{i-1} = U)\mathbb{1}_{(\mathbf{u}_1, \mathbf{u}_M)=y_{i-1}}\mathbb{P}_k(X_{i-1} = U|y_1^{i-1})}{\mathbb{P}_k(X_i = V|y_1^{i-1})} \cdot \\
 &\quad \mathbb{P}_k(X_i = V|y)\mathbb{1}_{(\mathbf{v}_1, \mathbf{v}_M)=y_i}, \\
 &= \frac{\mathbb{P}_k(X_{i-1} = U|y_1^{i-1})\mathbb{P}_k(X_i = V|y)}{\mathbb{P}_k(X_i = V|y_1^{i-1})} \cdot \\
 &\quad C^{(k)}(U, V)\mathbb{1}_{(\mathbf{u}_1, \mathbf{u}_M)=y_{i-1}}\mathbb{1}_{(\mathbf{v}_1, \mathbf{v}_M)=y_i}.
 \end{aligned}$$

## Equations de lissage 2 (EL2)

$$\begin{aligned}\mathbb{P}_k(X_{i-1} = U|y) &= \sum_{\mathbf{v} \in \mathcal{A}^M} \mathbb{P}_k(X_{i-1} = U, X_i = V|y), \\ &= \mathbb{P}_k(X_{i-1} = U|y_1^{i-1}) \mathbb{1}_{(\mathbf{u}_1, \mathbf{u}_M) = \mathbf{y}_{i-1}} \cdot \\ &\quad \sum_{\mathbf{v} \in \mathcal{A}^M} \frac{\mathbb{P}_k(X_i = V|y)}{\mathbb{P}_k(X_i = V|y_1^{i-1})} C^{(k)}(U, V) \mathbb{1}_{(\mathbf{v}_1, \mathbf{v}_M) = \mathbf{y}_i}.\end{aligned}$$