

Modèles d'arbre de Markov cachés

Yann Guédon⁽¹⁾ & Jean-Baptiste Durand⁽²⁾

(1) UMR CIRAD/CNRS/INRA/IRD/Université Montpellier II

Botanique et Bioinformatique de l'Architecture des Plantes

guedon@cirad.fr

(2) Laboratoire de Modélisation et Calcul

Institut d'Informatique et Mathématiques Appliquées de Grenoble

Jean-Baptiste.Durand@imag.fr

Contexte

Arborescence : structure orientée non-symétrique,

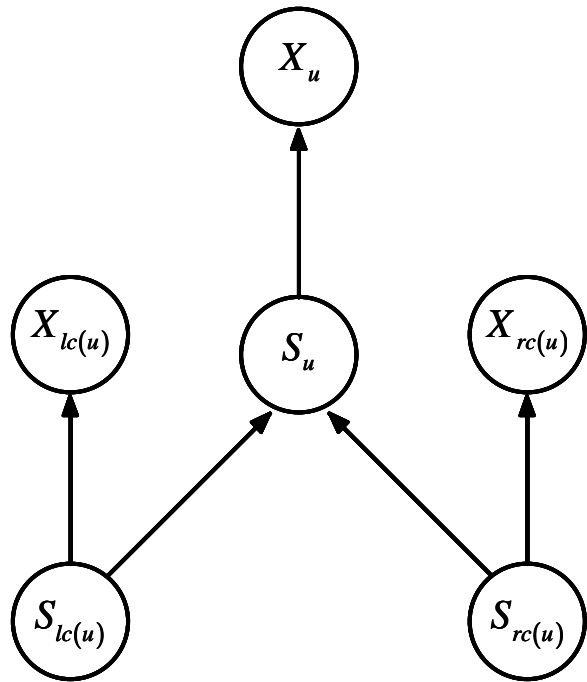
2 cas :

- arborescence orientée depuis **la racine**,
- arborescence orientée depuis **les feuilles**.

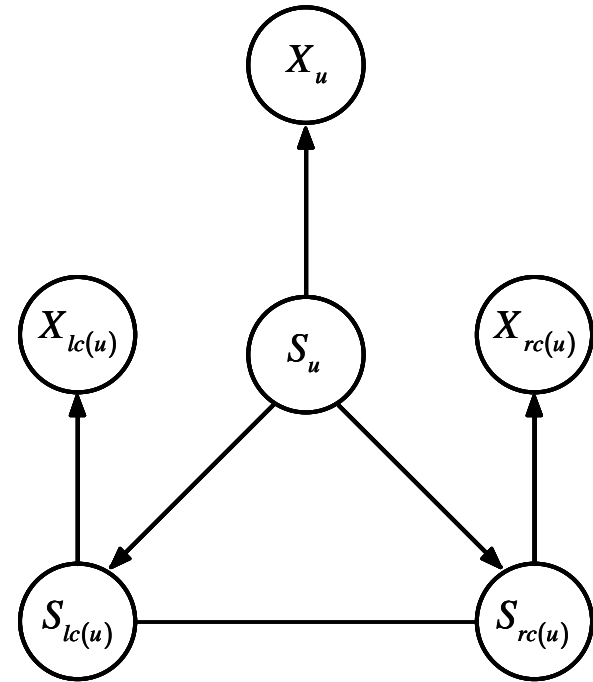
Topologie de l'arborescence fixée par les données.

Objectif : identification de zones homogènes (de nature arborescente), de ruptures dans des arborescences.

(a) hidden Markov in-tree model



(b) hidden Markov out-tree model



Définition d'un modèle d'arbre de Markov orienté depuis les feuilles (Markov in-tree model)

Arborescence binaire où $lc(u)$ et $rc(u)$ désignent les deux vertex fils du vertex u et $b(u)$ le vertex frère de u .

Un modèle d'arbre de Markov binaire, orienté depuis les feuilles à J états est défini par

- probabilités initiales (pour les feuilles) $\pi_j = P(S_u = j)$,
- probabilités de transition $p_{hij} = P(S_u = j | S_{lc(u)} = h, S_{rc(u)} = i)$.

J^2 lois de transition et $J^2 (J - 1)$ probabilités de transition indépendantes comme pour une chaîne de Markov d'ordre 2.

Dans le cas où l'arborescence est non-ordonnée, les probabilités de transition sont invariantes par permutation des vertex fils

$$\begin{aligned} p_{hij} &= P(S_u = j | S_{lc(u)} = h, S_{rc(u)} = i) \\ &= P(S_u = j | S_{lc(u)} = i, S_{rc(u)} = h) \\ &= p_{ihj}. \end{aligned}$$

Définition d'un modèle d'arbre de Markov orienté depuis la racine (Markov out-tree model)

Nous supposons que

$$\begin{aligned} & P \left(S_{lc(u)} = k, S_{rc(u)} = m | S_u = j \right) \\ \neq & P \left(S_{lc(u)} = k | S_u = j \right) P \left(S_{rc(u)} = m | S_u = j \right). \end{aligned}$$

Un modèle d'arbre de Markov binaire, orienté depuis la racine à J états est défini par

- probabilités initiales (pour la racine) $\pi_j = P(S_0 = j)$,
- probabilités de transition $p_{jkm} = P(S_{lc(u)} = k, S_{rc(u)} = m | S_u = j)$.

J lois de transition et $J(J - 1)$ probabilités de transition indépendantes.

Dans le cas où l'arborescence est non-ordonnée, les probabilités de transition sont invariantes par permutation des vertex fils

$$\begin{aligned} p_{jkm} &= P(S_{lc(u)} = k, S_{rc(u)} = m | S_u = j) \\ &= P(S_{lc(u)} = m, S_{rc(u)} = k | S_u = j) \\ &= p_{jmk}. \end{aligned}$$

Définition d'un modèle d'arbre de Markov orienté depuis la racine (Markov out-tree model)

Nous supposons que

$$\begin{aligned} & P \left(S_{lc(u)} = k, S_{rc(u)} = m | S_u = j \right) \\ &= P \left(S_{lc(u)} = k | S_u = j \right) P \left(S_{rc(u)} = m | S_u = j \right). \end{aligned}$$

Un modèle d'arbre de Markov, orienté depuis la racine à J états est défini par

- probabilités initiales (pour la racine) $\pi_j = P(S_0 = j)$,
- probabilités de transition $p_{jk} = P(S_{c(u)} = k | S_u = j)$.

J lois de transition et $J(J - 1)$ probabilités de transition indépendantes comme pour une chaîne de Markov d'ordre 1.

Définition d'un modèle d'arbre de Markov caché

Le processus d'observation $\{X_u\}$ discret univarié est relié au modèle d'arbre de Markov caché $\{S_u\}$ par les probabilités d'observation (ou d'émission)

$$b_j(y) = P(X_u = y | S_u = j)$$

Extension directe au cas multivarié incluant éventuellement des processus d'observation continus (avec par exemple des lois d'observation gaussiennes).

Application de l'algorithme EM

→ généralisation directe du cas des chaînes de Markov cachées.

Espérance conditionnelle de la log-vraisemblance des données complètes :

$$Q(\theta|\theta^{(k)}) = E \left\{ \log f \left(\bar{S}_0, \bar{X}_0; \theta \right) \mid \bar{X}_0 = \bar{x}_0; \theta^{(k)} \right\} .$$

Modèle d'arbre de Markov caché orienté depuis les feuilles

$$Q(\theta|\theta^{(k)}) = Q_\pi \left(\left\{ \pi_j \right\}_{j=0}^{J-1} \mid \theta^{(k)} \right) + \sum_{h,i} Q_p \left(\left\{ p_{hij} \right\}_{j=0}^{J-1} \mid \theta^{(k)} \right) \\ + \sum_j Q_b \left(\left\{ b_j(y) \right\}_{y=0}^{Y-1} \mid \theta^{(k)} \right).$$

avec

$$Q_\pi \left(\left\{ \pi_j \right\}_{j=0}^{J-1} \mid \theta^{(k)} \right) = \sum_j \sum_{u \in \text{feuilles}} P \left(S_u = j \mid \bar{X}_0 = \bar{x}_0; \theta^{(k)} \right) \log \pi_j,$$

$$Q_p \left(\left\{ p_{hij} \right\}_{j=0}^{J-1} \mid \theta^{(k)} \right) \\ = \sum_j \sum_{u \in \text{vertex internes}} P \left(S_u = j, S_{lc(u)} = h, S_{rc(u)} = i \mid \bar{X}_0 = \bar{x}_0; \theta^{(k)} \right) \log p_{hij}.$$

Modèle d'arbre de Markov caché orienté depuis la racine

$$Q(\theta|\theta^{(k)}) = Q_\pi \left(\{\pi_j\}_{j=0}^{J-1} | \theta^{(k)} \right) + \sum_j Q_p \left(\{p_{jkm}\}_{k,m=0}^{J-1} | \theta^{(k)} \right) + \sum_j Q_b \left(\{b_j(y)\}_{y=0}^{Y-1} | \theta^{(k)} \right).$$

avec

$$Q_\pi \left(\{\pi_j\}_{j=0}^{J-1} | \theta^{(k)} \right) = \sum_j P \left(S_0 = j | \bar{X}_0 = \bar{x}_0; \theta^{(k)} \right) \log \pi_j,$$

$$Q_p \left(\{p_{jkm}\}_{k,m=0}^{J-1} | \theta^{(k)} \right) = \sum_{k,m} \sum_{u \in \text{vertex internes}} P \left(S_{lc(u)} = k, S_{rc(u)} = m, S_u = j | \bar{X}_0 = \bar{x}_0; \theta^{(k)} \right) \log p_{jkm}$$

Dans les deux cas, pour les probabilités d'observation

$$Q_b \left(\{b_j(y)\}_{y=0}^{Y-1} \mid \theta^{(k)} \right) = \sum_y \sum_u P \left(X_u = y, S_u = j \mid \bar{X}_0 = \bar{x}_0; \theta^{(k)} \right) \log b_j(y).$$

Algorithme “montant-descendant” pour un modèle d'arbre de Markov caché orienté depuis les feuilles

Principe

$$\begin{aligned} & P(S_u = j | \bar{X}_0 = \bar{x}_0) \\ = & \frac{P(\bar{X}_{0 \setminus u} = \bar{x}_{0 \setminus u} | S_u = j)}{P(\bar{X}_{0 \setminus u} = \bar{x}_{0 \setminus u} | \bar{X}_u = \bar{x}_u)} P(S_u = j | \bar{X}_u = \bar{x}_u) \\ = & \alpha_u(j) \beta_u(j). \end{aligned}$$

Récurrance montante - Initialisation (feuilles de l'arborescence)

$$\begin{aligned}\beta_u(j) &= P(S_u = j | X_u = x_u) \\ &= \frac{b_j(x_u) \pi_j}{N_u}.\end{aligned}$$

Pour tous les vertex internes pris dans le sens montant,

$$\begin{aligned}\beta_u(j) &= P(S_u = j | \bar{X}_u = \bar{x}_u) \\ &= \frac{b_j(x_u) \sum_{h,i} p_{hij} \beta_{lc(u)}(h) \beta_{rc(u)}(i)}{N_u}.\end{aligned}$$

Le facteur de normalisation N_u est donné par

$$\begin{aligned} N_u &= P(X_u = x_u) \\ &= \sum_j b_j(x_u) \pi_j \end{aligned}$$

pour les feuilles, et

$$\begin{aligned} N_u &= P(X_u = x_u | \bar{X}_{c(u)} = \bar{x}_{c(u)}) \\ &= \sum_j b_j(x_u) \sum_{h,i} p_{hij} \beta_{lc(u)}(h) \beta_{rc(u)}(i) \end{aligned}$$

pour les vertex internes.

Récurrance descendante - Initialisation (racine de l'arborescence)

$$\alpha_0(j) = 1.$$

Pour tous les vertex restants pris dans le sens descendant,

$$\begin{aligned}\alpha_u(j) &= \frac{P(\bar{X}_{0 \setminus u} = \bar{x}_{0 \setminus u} | S_u = j)}{P(\bar{X}_{0 \setminus u} = \bar{x}_{0 \setminus u} | \bar{X}_u = \bar{x}_u)} \\ &= \frac{\sum_m \alpha_{\rho(u)}(m) b_m(x_{\rho(u)}) \sum_k p_{jkm} \beta_{b(u)}(k)}{N_{\rho(u)}}.\end{aligned}$$

Algorithme “montant-descendant” pour un modèle d'arbre de Markov caché orienté depuis la racine

Difficulté : récurrence montante dans le sens inverse de la structure.

→ Nécessité de calculer dans une passe préalable descendante les lois marginales $(P(S_u = j); j = 0, \dots, J - 1)$ pour chaque vertex u .

Récurrance montante - Initialisation (feuilles de l'arborescence)

$$\begin{aligned}\beta_u(j) &= P(S_u = j | X_u = x_u) \\ &= \frac{b_j(x_u) P(S_u = j)}{N_u}.\end{aligned}$$

Pour tous les vertex internes pris dans le sens montant,

$$\begin{aligned}\beta_u(j) &= P(S_u = j | \bar{X}_u = \bar{x}_u) \\ &= \left\{ \sum_{k,m} \frac{\beta_{lc(u)}(k) \beta_{rc(u)}(m) p_{jkm}}{P(S_{lc(u)} = k) P(S_{rc(u)} = m)} \right\} \frac{b_j(x_u) P(S_u = j)}{N_u}.\end{aligned}$$

Le facteur de normalisation N_u est donné par

$$\begin{aligned} N_u &= P(X_u = x_u) \\ &= \sum_j b_j(x_u) P(S_u = j) \end{aligned}$$

pour les feuilles, et

$$\begin{aligned} N_u &= \frac{P(\bar{X}_u = \bar{x}_u)}{P(\bar{X}_{lc(u)} = \bar{x}_{lc(u)}) P(\bar{X}_{rc(u)} = \bar{x}_{rc(u)})} \\ &= \sum_j \left\{ \sum_{k,m} \frac{\beta_{lc(u)}(k) \beta_{rc(u)}(m) p_{jkm}}{P(S_{lc(u)} = k) P(S_{rc(u)} = m)} \right\} b_j(x_u) P(S_u = j) \end{aligned}$$

pour les vertex internes.

Récurrance descendante - Initialisation (racine de l'arborescence)

$$\alpha_0(j) = 1.$$

Pour tous les vertex restants pris dans le sens descendant,

$$\begin{aligned}\alpha_u(j) &= \frac{P(\bar{X}_{0 \setminus u} = \bar{x}_{0 \setminus u} | S_u = j)}{P(\bar{X}_{0 \setminus u} = \bar{x}_{0 \setminus u} | \bar{X}_u = \bar{x}_u)} \\ &= \frac{1}{N_{\rho(u)} P(S_u = j)} \sum_k \frac{\beta_{b(u)}(k)}{P(S_{b(u)} = k)} \\ &\quad \times \sum_i p_{ijk} b_i(x_{\rho(u)}) \alpha_{\rho(u)}(i) P(S_{\rho(u)} = i).\end{aligned}$$

Cas des états des vertex fils conditionnellement indépendants (Durand *et al.*, 2004)

Récurrance montante

$$\begin{aligned}\beta_u(j) &= \left\{ \sum_{k,m} \frac{\beta_{lc(u)}(k) \beta_{rc(u)}(m) p_{jkm}}{P(S_{lc(u)}=k) P(S_{rc(u)}=m)} \right\} \frac{b_j(x_u) P(S_u=j)}{N_u} \\ &= \left\{ \sum_k \frac{\beta_{lc(u)}(k) p_{jk}}{P(S_{lc(u)}=k)} \right\} \left\{ \sum_m \frac{\beta_{rc(u)}(m) p_{jm}}{P(S_{rc(u)}=m)} \right\} \frac{b_j(x_u) P(S_u=j)}{N_u} \\ &= \frac{\beta_{u,lc(u)}(j) \beta_{u,rc(u)}(j) b_j(x_u) P(S_u=j)}{N_u}.\end{aligned}$$

Récurrance descendante

$$\begin{aligned}
 \alpha_u(j) &= \frac{1}{N_{\rho(u)} P(S_u = j)} \sum_k \frac{\beta_{b(u)}(k)}{P(S_{b(u)} = k)} \\
 &\quad \times \sum_i p_{ijk} b_i(x_{\rho(u)}) \alpha_{\rho(u)}(i) P(S_{\rho(u)} = i) \\
 &= \frac{1}{N_{\rho(u)} P(S_u = j)} \sum_i \left\{ \underbrace{\sum_k \frac{\beta_{b(u)}(k) p_{ik}}{P(S_{b(u)} = k)}}_{\beta_{\rho(u), b(u)}(i)} \right\} \\
 &\quad \times p_{ij} b_i(x_{\rho(u)}) \alpha_{\rho(u)}(i) P(S_{\rho(u)} = i)
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{P(S_u = j)} \sum_i p_{ij} \left\{ \underbrace{\frac{\beta_{\rho(u),u}(i) \beta_{\rho(u),b(u)}(i) b_i(x_{\rho(u)}) P(S_{\rho(u)} = i)}{N_{\rho(u)}}}_{\beta_{\rho(u)}(i)} \right\} \\
&\quad \times \frac{\alpha_{\rho(u)}(i)}{\beta_{\rho(u),u}(i)} \\
&= \frac{1}{P(S_u = j)} \sum_i \frac{p_{ij} \beta_{\rho(u)}(i) \alpha_{\rho(u)}(i)}{\beta_{\rho(u),u}(i)}.
\end{aligned}$$

Complexité en $O(J^2 n)$ au lieu de $O(J^3 n)$ dans le cas dépendant.

Algorithme “arrière-avant” pour une chaîne de Markov cachée

Principe

$$\begin{aligned} & P(S_t = j | X_0^{\tau-1} = x_0^{\tau-1}) \\ = & \frac{P(X_0^{t-1} = x_0^{t-1} | S_t = j)}{P(X_0^{t-1} = x_0^{t-1} | X_t^{\tau-1} = x_t^{\tau-1})} P(S_t = j | X_t^{\tau-1} = x_t^{\tau-1}) \\ = & \alpha_t(j) \beta_t(j). \end{aligned}$$

Récurrance arrière - $t = \tau - 1$:

$$\begin{aligned}\beta_{\tau-1}(j) &= P(S_{\tau-1} = j | X_{\tau-1} = x_{\tau-1}) \\ &= \frac{b_j(x_{\tau-1}) P(S_{\tau-1} = j)}{N_{\tau-1}}.\end{aligned}$$

$t = \tau - 2, \dots, 0$:

$$\begin{aligned}\beta_t(j) &= P(S_t = j | X_t^{\tau-1} = x_t^{\tau-1}) \\ &= \left\{ \sum_k \frac{\beta_{t+1}(k) p_{jk}}{P(S_{t+1} = k)} \right\} \frac{b_j(x_t) P(S_t = j)}{N_t}.\end{aligned}$$

Le facteur de normalisation N_t est donné par

$$\begin{aligned} N_{\tau-1} &= P(X_{\tau-1} = x_{\tau-1}) \\ &= \sum_j b_j(x_{\tau-1}) P(S_{\tau-1} = j) \end{aligned}$$

pour $t = \tau - 1$

$$\begin{aligned} N_t &= \frac{P(X_t^{\tau-1} = x_t^{\tau-1})}{P(X_{t+1}^{\tau-1} = x_{t+1}^{\tau-1})} \\ &= \sum_j \left\{ \sum_k \frac{\beta_{t+1}(k) p_{jk}}{P(S_{t+1} = k)} \right\} b_j(x_t) P(S_t = j) \end{aligned}$$

pour $t < \tau - 1$.

Récurrance avant - $t = 0$:

$$\alpha_0(j) = 1.$$

$t = 1, \dots, \tau - 1$:

$$\begin{aligned}\alpha_t(j) &= \frac{P(X_0^{t-1} = x_0^{t-1} | S_t = j)}{P(X_0^{t-1} = x_0^{t-1} | X_t^{\tau-1} = x_t^{\tau-1})} \\ &= \frac{1}{N_{t-1} P(S_t = j)} \sum_i p_{ij} b_i(x_{t-1}) \alpha_i(t-1) P(S_{t-1} = i).\end{aligned}$$

Remarques :

- Tous les vertex de la sous-arborescence \bar{x}_u sont visités avant sa racine x_u dans la récurrence montante,
- Certains des vertex de $\bar{x}_{0 \setminus u}$ n'ont pas encore été visités dans la récurrence descendante à l'étape u .

Dans la première passe obligatoire (“avant” ou “arrière” dans le cas de séquences, “montante” dans le cas d'arborescences) - qui peut être précédée par une passe préalable de calcul des lois marginales suivant la direction de la structure - il est seulement possible de fusionner de l'information.

→ Ceci implique que le demi-degré intérieur (nombre d'arcs entrants) ou extérieur (nombre d'arcs sortants) de chaque vertex est égal à 1.

Pour l'application à l'analyse de la structure des plantes

Objectif : limiter le nombre de probabilités de transition indépendantes.

- Vertex fils non-ordonnés ou partiellement ordonnés (vertex fils correspondant à l'entité suivante - construite par le même méristème - distingué des autres vertex fils correspondant aux entités portées),
- Succession d'états transitoires et état final absorbant (variable d'état ordinale); ordre de grandeur du nombre d'états : 5,
- Prise en compte uniquement des entités jouant un rôle dans la construction de la structure de la plante (par opposition aux entités jouant principalement un rôle d'assimilation),

- Combinaison des probabilités de transition estimées du type $p_{hij} = P(S_u = j | S_{c(u)} = \{h, i\})$ et $p_{ghij} = P(S_u = j | S_{c(u)} = \{g, h, i\})$ dans le cas de modèles d'arbre de Markov cachés orientés depuis les feuilles, ou $p_{jkm} = P(S_{c(u)} = \{k, m\} | S_u = j)$ et $p_{jkmn} = P(S_{c(u)} = \{k, m, n\} | S_u = j)$ dans le cas de modèles d'arbre de Markov cachés orientés depuis la racine.

Les modèles d'arbre de Markov cachés orientés depuis les feuilles permettent de prendre en compte des arborescences où seules les parties périphériques de la plante sont observables du fait de l'effacement des marqueurs morphologiques lié à l'élagage et à la croissance cambiale (en diamètre).