

INFÉRENCE STATISTIQUE SUR DES PROCESSUS DE BRANCHEMENTS MULTI-TYPES

Amel Ouaari ¹ & Rachid Senoussi ²

¹ *Biostatistique et Processus Spatiaux (BioSP), INRA, Site Agroparc, 84914 Avignon, France. et amel.ouaari@paca.inra.fr*

² *Biostatistique et Processus Spatiaux (BioSP), INRA, Site Agroparc, 84914 Avignon, France. et rachid.senoussi@paca.inra.fr*

Résumé.

On s'intéresse aux estimateurs du maximum de vraisemblance des paramètres associés à des processus de branchements en temps continu, multi-types et homogènes. L'exposé passe par la description probabiliste et paramétrée de l'évolution temporelle de plusieurs populations ou classes en interaction. Nous décrivons le comportement du système multi-types via l'évolution des fonctions génératrices de la loi des effectifs et faisons le lien avec la théorie des équations différentielles aux dérivées partielles. Notre approche statistique consiste à inverser théoriquement ou numériquement la fonction génératrice en utilisant le théorème de Cauchy pour les fondrées analytiques pour inférer ces paramètres.

Mots-clés. Processus de branchements en temps continu multi-types, estimateurs de maximum de vraisemblance, équations aux dérivées partielles, théorème de Cauchy.

Abstract.

This paper deals with the maximum-likelihood estimators of parameters associated to homogeneous multi-type continuous-time branching processes. It describes, in a probabilistic way, the temporal change of sever interacting populations or classes. We first describe the behaviour of the multi-type system through the evolution of the generating functions of the offspring distribution and link it to Partial Differential Equation theory. Our statistical approach consists in inverting theoretically or numerically the probability-generating functions using Cauchy's theorem in order to infer these parameters.

Keywords. Multi-type continuous-time branching processes, maximum-likelihood estimators, Partial Differential Equations, Cauchy's theorem.

1 Introduction

L'intérêt des processus de branchement multi-types réside dans les nombreuses applications dans des domaines aussi divers que la biologie, l'épidémiologie, la démographie, la physique des particules, etc. Les processus de branchements sont considérés comme des modèles de probabilité appropriés pour la description du comportement des systèmes dont

les composantes (cellules, particules, \dots etc) se reproduisent, se transforment, interagissent ou meurent. Par ailleurs, en dynamique de populations, il y a un manque flagrant de lois de probabilité modélisant des comptages simultanés des individus de plusieurs populations interactives et évoluant dans le temps. L'intérêt de développer des méthodes d'inférence statistique propres à ces modèles est de grande actualité. Il existe de nombreuses approches d'estimation paramétriques dans les processus de branchements, mais celles-ci sont pour la plupart basées sur des méthodes de contrastes et de propriétés de martingale [6]. Nous proposons ici de revenir à la méthode du maximum de vraisemblance en inversant les équations des fonctions génératrices.

Notations:

Nous utilisons les notations suivantes:

K sera le nombre de populations en jeu, et $\mathbf{i} = (i_1, \dots, i_K)$, multi-indice de \mathbb{N}^K , décrira les tailles respectives des différents sous populations.

De même, $\mathbf{0} = (0, \dots, 0)$, $\mathbf{1} = (1, \dots, 1)$ et $\mathbf{e}_i = (0, \dots, 0, \overbrace{1}^i, 0, \dots, 0)$: seront des vecteurs de \mathbb{N}^K d'utilisation constante.

On notera $\mathbf{X}(t) = (X_1(t), \dots, X_K(t)) \in \mathbb{N}^K$ les tailles des populations à l'instant t et si $\mathbf{z} = (z_1, \dots, z_K)$ est dans \mathbb{C}^K , on écrira $\mathbf{z}^{\mathbf{i}} = z_1^{i_1} \times \dots \times z_K^{i_K}$

2 Equations du branchement multi-types

On considère les K fonctions génératrices de $X(t)$, quand on part avec un état initial $X(0) = \mathbf{e}_i$, $i = 1, \dots, K$

$$G_i(t, \mathbf{z}) = E(\mathbf{z}^{\mathbf{X}(t)} / X(0) = \mathbf{e}_i).$$

Si on ne considère que des branchements homogènes dans le temps, alors la notation vectorielle: $G(t, \mathbf{z}) = (G_1(t, \mathbf{z}), \dots, G_K(t, \mathbf{z}))$ permet de traduire simplement l'équation de Chapman-Kolmogorov du branchement multi-type par:

$$G(t + s, \mathbf{z}) = G(t, G(s, \mathbf{z})), \quad \text{pour tous } s, t \geq 0 \text{ et } \mathbf{z} \in D(0, 1)^K \subset \mathbb{C}^K$$

avec la condition initiale $G(0, \mathbf{z}) = \mathbf{z}$.

Remarques:

- Si $X(0) = \mathbf{m} = (m_1, \dots, m_K)$ alors $G_{\mathbf{m}}(t, \mathbf{z}) = G_1^{m_1}(t, \mathbf{z}) \times \dots \times G_K^{m_K}(t, \mathbf{z})$.
- Si $X(0)$ est de loi γ_0 sur \mathbb{N}^K et de génératrice G_0 , alors: $G(t, \mathbf{z}) = G_0(G(t, \mathbf{z}))$.

$X(t)$ est un processus markovien de sauts, homogène dans le temps. On peut de manière explicite exprimer le générateur infinitésimal Q du semi groupe

$$P_{i,j}(t) = P(X(t) = j \mid X(0) = i)$$

$$Q_{i,j} = \lim_{t \rightarrow 0} \frac{P_{i,j}(t) - \mathbb{I}_{\{i=j\}}}{t}$$

On a alors pour $\mathbf{z} \in D(0,1)^K$ et pour $V_j(\mathbf{z}) = \partial_t G_j(0, \mathbf{z})$, l'égalité suivante:

$$V_j(\mathbf{z}) = Q_{j,j} \mathbf{z}^j + \sum_{l \neq j} Q_{j,l} \mathbf{z}^l.$$

En écriture vectorielle, on en déduit pour $V(\mathbf{z}) = \partial_t G(0, \mathbf{z})$, que

$$\partial_t G(t, \mathbf{z}) = -V(\mathbf{z})(\partial_z \mathbf{G}(t, \mathbf{z}))^T. \quad (1)$$

Notons que, chaque colonne de l'équation dynamique (en t) ne met en jeu qu'un seul G_i à la fois (et les dérivées en $t = 0$ des autres G_i).

Il s'agit donc de K équations aux dérivées partielles, linéaires homogènes, identiques, dont, seules les conditions initiales varient.

Fonction génératrice des moments

Rappelons que si X est une v.a sur \mathbb{N}^K de loi $P = (p_{\mathbf{m}}, \mathbf{m} \in \mathbb{N}^K)$ a pour génératrice

$$G(\mathbf{z}) = E(\mathbf{z}^{\mathbf{X}(t)}) = \sum p_{\mathbf{m}} \mathbf{z}^{\mathbf{m}}$$

alors, G est (multi) analytique en \mathbf{z} dans $(D(0,1))^K$, et que les dérivés partielles d'ordre \mathbf{m} en 0, redonnent la loi de probabilité P :

$$\partial_{z_1^{m_1} \dots z_K^{m_K}} G(0) = m_1! \dots m_K! p_{\mathbf{m}}.$$

3 Approche statistique par Max-Vraisemblance

Pour l'interprétation du modèle de branchement, la paramétrisation des lois de transition $P_{i,j}(t)$ ou de façon équivalente des génératrices $G_i(t, \mathbf{z})$ se fera naturellement par le biais du générateur infinitésimal $Q_{i,j}$, autrement dit des fonctions analytiques $V_j(\mathbf{z})$.

- $-Q_{j,j}$: est le paramètre de la loi exponentielle du temps de vie d'un individu de la population j .
- $\frac{Q_{j,\mathbf{m}}}{-Q_{j,j}}$: est la probabilité pour qu'un individu de la population j produise $\mathbf{m} = (m_1, \dots, m_K)$ descendants lors de son décès .

Le problème d'inférence réside donc dans l'inversion des fonctions génératrices après avoir résolu les équations dynamiques (1) afin d'extraire les probabilités correspondantes.

3.1 Vraisemblance avec des observations ponctuelles de $X(t)$ dans le temps t_1, \dots, t_m

La propriété de Markov dit que:

$$\mathcal{L}(X(t_1), \dots, X(t_m)/X(0)) = \mathcal{L}(X(t_1)/X(0)) \cdots \mathcal{L}(X(t_m)/X(t_{m-1}))$$

et l'hypothèse d'homogénéité nous dit aussi que pour les observations dans \mathbb{N}^K

$$X(0) = n_0, X(t_1) = n_1, \dots, X(t_m) = n_m.$$

La vraisemblance s'écrit :

$$V_\theta = P_{n_0, n_1}(t_1 | \theta) P_{n_1, n_2}(t_2 - t_1 | \theta) \cdots P_{n_{m-1}, n_m}(t_m - t_{m-1} | \theta)$$

Lorsque $\theta = (\theta_1, \dots, \theta_r)$ désigne le vecteur des paramètres d'intérêt de dimension finie r spécifiant le générateur infinitésimal.

L'estimateur du maximum de vraisemblance s'obtient par la résolution des r équations:

$$\partial_\theta \log V = \sum_{i=1}^m \frac{\partial_{\theta_i} P_{n_{i-1}, n_i}(t_i - t_{i-1} | \theta)}{P_{n_{i-1}, n_i}(t_i - t_{i-1} | \theta)} = 0.$$

Ce qui nous amène à calculer dans un premier temps les fonctions génératrices des transitions $P_{n, \cdot}(s | \theta)$ associées à:

$$G^n(s, \mathbf{z} | \theta) = G_1^{n_1}(s, \mathbf{z} | \theta) \times \cdots \times G_K^{n_K}(s, \mathbf{z} | \theta)$$

c'est-à-dire

$$P_{n, m}(s | \theta) = \frac{1}{m!} \partial_{z_1^{m_1} \dots z_K^{m_K}} G^n(s, \mathbf{z} | \theta)$$

3.2 Inversion de la génératrice

Si G est une fonction analytique, pour tout $\mathbf{z} \in D(0, r)^K$, l'utilisation du théorème de Cauchy pour les fonctions analytiques (uni ou multivariées) permet de calculer les probabilités correspondantes

$$P_{n_1 \dots n_K} = \frac{1}{(2\pi i)^K} \int_{\gamma_1} \int_{\gamma_2} \int_{\gamma_K} \frac{G(\xi_1, \dots, \xi_K)}{\xi_1^{n_1+1}, \dots, \xi_K^{n_K+1}} d\xi_1 \cdots d\xi_K$$

L'approximation sera d'autant meilleure quand K n'est pas très grand et les chemins γ_i simples.

3.3 Application

Partant de l'équation aux dérivées partielles:

$$\begin{cases} \frac{\partial G_i(t, \mathbf{z})}{\partial t} &= \sum_{j=1}^k -V_j(\mathbf{z}) \frac{\partial G_i(t, \mathbf{z})}{\partial z_j}, \\ G_i(0, \mathbf{z}) &= z_i. \end{cases}$$

On normalisera d'abord les écritures des équations car pour tout Q^* vérifiant

$$Q^* \geq \sup(\bar{Q}_{jj}) \quad \text{où} \quad \bar{Q}_{jj} = - \sum_{\substack{r=0 \\ r \neq j}}^{\infty} Q_{jr}, \quad (j = 1, \dots, K),$$

et pour tout j , il existe $Q_{j,j}^{**}$ tel que $\bar{Q}_{jj} + Q_{j,j}^{**} = Q^*$ et alors:
 $V_j(\mathbf{z}) = Q^*[H_j(\mathbf{z}) - z_j]$ avec $H_j(\mathbf{z})$ désignant la fonction génératrice du branchement de type j , mieux interprétable.

Exemple de branchements Poissonniens indépendants

$$H_j(\mathbf{z}) = \prod_{i=1}^K e^{\lambda_{ji}(z_j-1)} = e^{\sum_{i=1}^K \lambda_{ji} z_j} e^{-\sum_{i=1}^K \lambda_{ji}} = e^{\langle \bar{\lambda}_j, \mathbf{z} \rangle} e^{-\bar{\lambda}_j}$$

où $\bar{\lambda}_j = \sum_{i=1}^K \lambda_{ji}$

On standardisera l'EDP de la façon suivante:

Soit $\Lambda = (\lambda_{ij})$ la matrice des paramètres,

et $\Lambda^{-1} = (\mu_{ij})$, $\gamma_j = e^{-\bar{\lambda}_j}$, $\beta_{jl} = \sum_{j=1}^k \sum_{l=1}^k \lambda_{jl} \mu_{lj}$

alors $H_j(\Lambda \mathbf{u}) = e^{u_j - \bar{\lambda}_j}$

Le changement de variables $\mathbf{z} = \Lambda \mathbf{u}$ donne la fonction $\tilde{G}(t, \mathbf{u}) = G(t, \Lambda \mathbf{u})$.

Elle vérifie: $\partial_t \tilde{G}(t, \mathbf{u}) = \theta \sum_{l=1}^k V_j(\mathbf{u}) \partial_{u_l} \tilde{G}(t, \mathbf{u})$

avec

$$V_j(\mathbf{u}) = \gamma_{j1} e^{u_1} + \dots + \gamma_{jj} e^{u_K} - (\beta_{j1} u_1 + \dots + \beta_{jK} u_K).$$

La résolution du système d'équations différentielles ordinaires de dimension K :

$$\begin{cases} u'_1 = V_1(u_1, \dots, u_K) \\ \vdots \\ u'_K = V_K(u_1, \dots, u_K) \end{cases}$$

s'écrit $\mathbf{u}' = \gamma e^{\mathbf{u}} - \beta \mathbf{u}$
avec $\mathbf{u} = (u_1, \dots, u_k)^T$, $e^{\mathbf{u}} = (e^{u_1}, \dots, e^{u_k})$, $\beta = \beta_{jl}$

Cas d'une population unique

Pour résoudre l'EDP (1), il faut alors résoudre l'équation différentielle ordinaire:

$$\frac{du}{dt} = \gamma e^u - \beta u$$

c'est-à-dire chercher une primitive F de $f(u) = \frac{1}{\gamma e^u - \beta u}$ et une intégrale première de ce système.

La génératrice solution sera la valeur z_0 satisfaisant l'équation implicite

$$F(z) - F(z_0) = t$$

Nous concluons la présentation par des calculs numériques et des simulations illustrant la portée de cette approche d'estimation.

Bibliographie

- [1] Athreya . P. E. Ney (1972), Branching Processes, Springer, Berlin.
- [2] Athreya, K. B. et Niels Keiding(1977), Estimation theory for continuous-time branching processes, Indian Institute of Science, Bangalore, India, University of Copenhagen, Denmark
- [3] González, C. et Minuesa, I. del Puerto(2004), Maximum likelihood estimation and expectation-maximization algorithm for controlled branching processes, Data Analysis, Badajoz, Spain.
- [4] González, M. et Martín, J. et Martínez, R. et Mota, M. (2007), Non-parametric Bayesian estimation For multitype branching processes through simulation-based methods, Computational Statistics Data Analysis, Badajoz, Spain.
- [5] Neils Becker (1977), Estimation for Discrete Time Branching Processes with Application to Epidemics , Biometrics, La Trobe university, Bundoora, Australia.
- [6] Soren Asemussem et Neils Keiding (1978), Martingale central limit theorems and asymptotic estimation theory for multitype branching processes, University of Copenhagen.