

# MODÈLES MULTI-ÉTATS MARKOVIENS ET SEMI-MARKOVIENS

**Philippe Saint Pierre**

Institut de Mathématiques de Toulouse

Université Toulouse III - Paul Sabatier

[philippe.saint-pierre@math.univ-toulouse.fr](mailto:philippe.saint-pierre@math.univ-toulouse.fr)

# ANALYSE DE SURVIE

## EXEMPLES

- Epidémiologie : temps avant l'apparition d'une maladie,
- Assurance : temps avant un sinistre,
- Fiabilité : temps avant la panne d'une machine.

## CONTEXTE

- On cherche à étudier une durée  $T$ .
- La variable  $T$  est souvent **censurée**
  - fin de l'étude,
  - perdus de vus.
- Censure à droite : on observe  $(Y_i, \delta_i)_{1 \leq i \leq n}$ ,

$$\begin{cases} Y_i &= \inf(T_i, C_i), \\ \delta_i &= \mathbf{1}_{T_i \leq C_i}. \end{cases}$$

## OBJECTIF

On cherche à estimer

$$S(t) = \mathbb{P}(T \geq t).$$

- Absence de censure : on observe  $Y_1, \dots, Y_n$ .

Estimateur empirique :  $\hat{S}_{emp}(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Y_i \geq t}$ .

Pas utilisable en présence de censure : **biais**

- **Hypothèse de censure indépendante** : la censure n'apporte aucune information sur l'événement étudié.
- **Hypothèse forte**, notamment pour les "perdus de vue"

Ex : Asthme, VIH, Etude Prozac, ...

## ESTIMATION NON PARAMÉTRIQUE : KAPLAN-MEIER

$$\hat{S}(t) = \prod_{Y_i \leq t} \left( 1 - \frac{n_i}{m_i} \right),$$

$n_i$  est le nombre de sujets à risque au temps  $Y_i$ ,

$m_i$  est le nombre de décès au temps  $Y_i$ .

## ESTIMATEUR DE KAPLAN-MEIER (BIS)

$$\hat{S}(t) = \frac{1}{n} \sum_{i=1}^n \hat{W}_i \mathbf{1}_{Y_i \geq t},$$

où  $\hat{W}_i = \frac{\delta_i}{\hat{G}(Y_i)}$ ,  $G(t) = \mathbb{P}(C \geq t)$  et  $\hat{G}$  son estimateur Kaplan-Meier.

## MODÈLE SEMI-PARAMÉTRIQUE DE COX

- Etudier l'effet des covariables sur la durée de survie.
- **Modèle à risques proportionnels**

$$\alpha(t | Z = Z_i) = \alpha_0(t) \exp(\beta^T Z_i)$$

$\beta$  vecteur de paramètres  $Z_i$  vecteur de covariable associé à l'individu  $i$

- Hypothèse de **proportionnalité des risques** ( $\beta$  indépendant du temps)

$$\frac{\alpha(t | Z = 1)}{\alpha(t | Z = 0)} = \exp(\beta)$$

- Interprétation en terme de **risque relatif** :  $\exp(\beta)$

## MODÈLE SEMI-PARAMÉTRIQUE DE COX

- Coefficients de régression  $\beta$  estimés par maximisation de la **vraisemblance partielle de Cox**

$$L_{Cox} = \prod_{i=1}^n \prod_{t \geq 0} \left( \frac{Y_i(t) \exp(\beta^T Z_i)}{\sum_{j=1}^n Y_j(t) \exp(\beta^T Z_j)} \right)^{\delta_i(t)}$$

- $Y_i(t) = 1$  si l'individu  $i$  est à risque au temps  $t$  et 0 sinon.
- $\delta_i(t) = 1$  si l'individu  $i$  subit l'événement au temps  $t$ .

## MODÈLES PARAMÉTRIQUES

- **Lois paramétriques** pour le taux de hasard (exponentiel, Weibull, ...)
- **Cox** parametric model

$$\alpha(t | Z) = \alpha_0(t) \exp(\beta^T Z) \quad \text{ou} \quad S(t | Z) = (S_0(t))^{\exp(\beta^T Z)}$$

- **Accelerated Failure Time** model (AFT)

$$\log(T) = \alpha + \beta^T Z + \varepsilon \quad \text{ou} \quad S(t | Z) = S_0(te^{\beta^T Z})$$

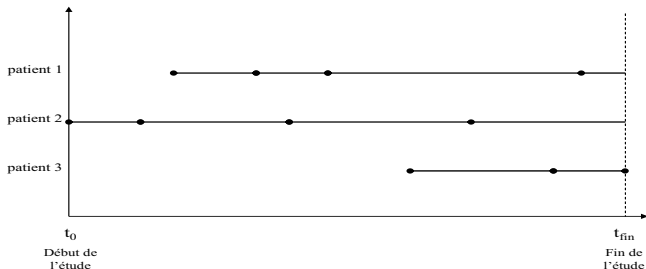
- Estimation par maximum de vraisemblance

$$L = \prod_{i=1}^n (S_i(t) \alpha_i(t))^{\delta_i} \times S_i(t)^{1-\delta_i}.$$

# MODÈLES MULTI-ÉTATS

## TYPE DE DONNÉES

- **Mesures répétées** dans le temps : la même variable est mesurée plusieurs fois pour un même individu : **données longitudinales**
- Mesures effectuées à des moments quelconques (spécifiques à chaque individu)





## MODÈLES MULTI-ÉTATS

- Les **modèles multi-états** permettent de modéliser ce type de données
- Notions d'état et de processus pour représenter l'évolution d'un phénomène
- L'analyse consiste à estimer les **intensités de transition** entre les états

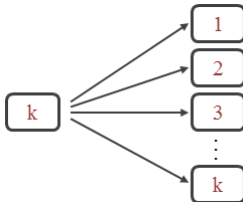
## Analyse de survie



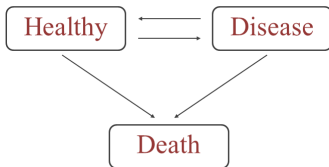
## Événements répétés



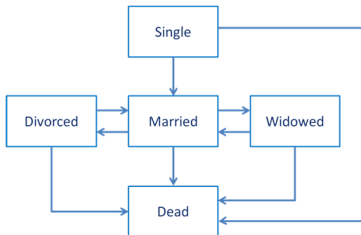
## Événements concurrents



Épidémiologie : représenter l'évolution d'un patient à travers les différents stades d'une maladie



Science sociale : représenter l'évolution de la situation professionnelle ou familiale



- **Propriété de Markov** : l'information sur les états précédents est résumée par l'état présent
- Les intensités de transition  $\alpha$  peuvent dépendre de deux échelles de temps :
  - la **durée du suivi**  $t$  (ou le temps calendaire, l'âge),
  - le **temps de séjour**  $d$

$\alpha(t, d) = \alpha$	→	Modèle de <b>Markov homogène</b>
$\alpha(t, d) = \alpha(t)$	→	Modèle de <b>Markov non-homogène</b>
$\alpha(t, d) = \alpha(d)$	→	Modèle <b>semi-Markovien homogène</b>
$\alpha(t, d) = \alpha(t, d)$	→	Modèle <b>semi-Markovien non-homogène</b>

# MODÈLE DE MARKOV

## PROCESSUS DE MARKOV

- $\{X(t), t \in [0, +\infty[$  un **processus de Markov** à temps continu et à espace d'états fini  $S = \{1, \dots, s\}$
- Probabilités de transition  $\mathbf{P}(s, t) = \{p_{hj}(s, t)\}$

$$p_{hj}(s, t) = \mathbb{P}(X(t) = j \mid X(s) = h)$$

- Intensités de transition entre les états,  $\mathbf{Q}(t) = \{\alpha_{hj}(t)\}$

$$\alpha_{hj}(t) = \lim_{\Delta t \rightarrow 0} \frac{p_{hj}(t, t + \Delta t) - p_{hj}(t, t)}{\Delta t}, \quad h \neq j$$

$$\alpha_{hh}(t) = - \sum_{j \neq h} \alpha_{hj}(t)$$

- Considérons  $X_1(\cdot), \dots, X_n(\cdot)$ ,  $n$  copies indépendantes de  $X(\cdot)$

## MODÈLE DE MARKOV HOMOGENÈ

- La matrice des probabilités de transition

$$\mathbf{P}(s, s + \Delta t) = \mathbf{P}(0, \Delta t) = \exp(\mathbf{Q} \times \Delta t)$$

- Intensités de transition :  $\alpha_{hj}(t) = \alpha_{hj}$
- Modèle à risques proportionnels

$$\alpha_{hj} = \alpha_{hj0} \exp(\beta_{hj}^T \mathbf{Z}), \quad h \neq j$$

$\alpha_{hj0}$  intensité de transition de base

$\mathbf{Z} = (Z_1, \dots, Z_k)^T$  vecteur de  $k$  covariables associé à l'individu

$\beta_{hj} = (\beta_{hj,1}, \dots, \beta_{hj,k})^T$  vecteur de coefficients de régression associé à la transition de l'état  $h$  vers l'état  $j$

## MODÈLE DE MARKOV HOMOGÈNE PAR PÉRIODE

- Intensités de transition définies sur  $r + 1$  intervalles  $[t_{k-1}, t_k[$

$$\begin{aligned} \alpha_{hj}(t) &= \alpha_{hj0} && \text{si } t_0 \leq t < t_1 \\ \alpha_{hj}(t) &= \alpha_{hjk} = \alpha_{hj0} \exp\left(\sum_{v=1}^k \gamma_{hj,v}\right) && \text{si } t_k \leq t < t_{k+1} \end{aligned}$$

$\forall k = 1, \dots, r$  et  $t_{r+1} = +\infty$

$\alpha_{hjk}$  intensités de base dans l'intervalle  $[t_k, t_{k+1}[$

$\gamma_{hj} = (\gamma_{hj,1}, \dots, \gamma_{hj,k})$  vecteur de coefficients de régression

- Les intensités sont **constantes** sur chaque intervalle
- Estimation des paramètres par maximum de vraisemblance

## MODÈLE DE MARKOV NON-HOMOGÈNE

- Intensités de transition  $\alpha_{hj}(t)$  dépendent de la durée du suivi
- $\mathbf{N} = (N_{hj}; h \neq j, h, j = 1, \dots, s)$  un processus de comptage multivarié

$N_{hj}(t)$  compte le nombre de transitions (dans toute la population) observées de  $h$  vers  $j$  dans  $[0, t]$

- $\lambda = (\lambda_{hj}; h \neq j)$ , intensité du processus  $\mathbf{N}$

$$\lambda_{hj}(t) = \alpha_{hj}(t)Y_h(t), \forall h \neq j$$

$Y_h(t) = \sum_{i=1}^n \mathbb{1}_{\{X_i(t^-)=h\}}$  est le nombre total d'individus observés dans  $h$  juste avant  $t$  (individus à risque)



## MODÈLE DE MARKOV NON-HOMOGÈNE

- $\mathbf{A}(t) = \{A_{hj}(t)\}$ , la matrice des intensités cumulées

$$A_{hj}(t) = \int_0^t \alpha_{hj}(s) ds$$

- La matrice des probabilités de transition  $\mathbf{P}(s, t) = \{p_{hj}(s, t)\}$  est définie par le produit intégral

$$\mathbf{P}(s, t) = \prod_{u \in ]s, t]} (\mathbf{Id} + d\mathbf{A}(u))$$

**Id** matrice identité

## ESTIMATION NON-PARAMÉTRIQUE

- Estimateur de **Nelson-Aalen** des intensités cumulées

$$\widehat{\mathbf{A}}_{hj}(t) = \int_0^t \frac{J_h(u)}{Y_h(u)} dN_{hj}(u), \quad \forall h \neq j$$

$$J_h(t) = \mathbf{1}_{\{Y_h(t) > 0\}}$$

- Estimateur de **Aalen-Johansen** de la matrice des probabilités de transition

$$\widehat{\mathbf{P}}(s, t) = \prod_{u \in ]s, t]} \left( \mathbf{Id} + d\widehat{\mathbf{A}}(u) \right), \quad 0 < s \leq t$$

$\widehat{\mathbf{A}}$  estimateur de Nelson-Aalen

- Comparaison des intensités obtenues dans différents groupes (test du Log-rank)

## ESTIMATION SEMI-PARAMÉTRIQUE

- Les intensités de transition suivent un modèle à risques proportionnels

$$\alpha_{hji}(t | \mathbf{Z}_i) = \alpha_{hj0}(t) \exp(\beta_{hj}^T \mathbf{Z}_i) \quad , h \neq j$$

- Sachant  $\beta$ , l'estimateur de Breslow

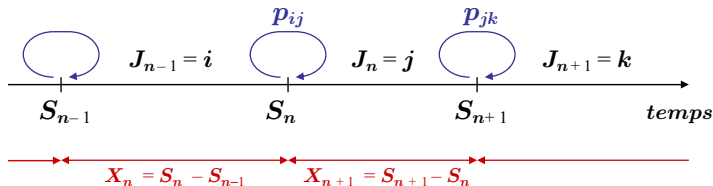
$$\hat{A}_{hj0}(t) = \int_0^t \frac{J_h(u)}{\sum_{i=1}^n \exp(\beta_{hj}^T \mathbf{Z}_i) Y_{hi}(t)} dN_{hj}(u),$$

- $\beta$  estimé par maximisation de la vraisemblance partielle de Cox

$$\mathcal{L}_{Cox}(\beta) = \prod_t \prod_{i=1} \prod_{h \neq j} \left[ \frac{Y_{hi}(t) \exp(\beta_{hj}^T \mathbf{Z}_i)}{\sum_{i=1}^n \exp(\beta_{hj}^T \mathbf{Z}_i) Y_{hi}(t)} \right]^{\Delta N_{hj}(t)}$$

# MODÈLE SEMI-MARKOVIENT HOMOGÈNE

- Modèle Markovien inadapté  $\implies$  Approche semi-Markovienne
- Soit  $(J_n, S_n)_{n \geq 0}$  un processus semi-Markovien



- $(J_n)_{n \geq 0}$  est une chaîne de Markov homogène
- $X_n = S_n - S_{n-1}$ , le temps de séjour dans l'état  $J_{n-1}$
- $J_n \neq J_{n+1}$

- Probabilités de transition de la chaîne de Markov  $(J_n)_{n \geq 0}$

$$p_{hj} = \mathbb{P}(J_{n+1} = j \mid J_n = h)$$

- Distribution du temps de séjour (dans l'état  $h$  avant d'aller dans l'état  $j$ )

$$F_{hj}(d) = \mathbb{P}(X_{n+1} \leq d \mid J_n = h, J_{n+1} = j)$$

$\implies S_{hj}(x), f_{hj}(x), \alpha_{hj}(x)$  (fonctions de survie, de densité et d'intensité correspondantes)

- Intensités de transition du processus semi-Markovien

$$\lambda_{hj}(d) = \lim_{\Delta d \rightarrow 0} \frac{1}{\Delta d} \mathbb{P}(J_{n+1} = j, d < X_{n+1} \leq d + \Delta d \mid J_n = h, X_{n+1} > d)$$

$$= \begin{cases} \frac{p_{hj}f_{hj}(d)}{S_h(d)} & \text{si } S_h(d) = \sum_{j=1}^s p_{hj}S_{hj}(d) > 0 \\ 0 & \text{sinon} \end{cases}$$

## VRAISEMBLANCE

- Contributions à la vraisemblance

- Transition de  $h \rightarrow j$  observée :  $S_h(d)\lambda_{hj}(d) = p_{hj}f_{hj}(d)$
- Censure à droite dans l'état  $h$  :  $S_h(d)$

- Vraisemblance

$$L = \prod_{i=1}^n \left\{ \prod_{k=1}^{n_i} p_{J_{i,k-1}J_{i,k}} f_{J_{i,k-1}J_{i,k}}(X_{i,k}) \right\} \times \left[ S_{J_{i,n_i}}(U_i) \right]^{\delta_i}$$

$\delta_i = 1$  si l'individu  $i$  est censuré et 0 sinon

$U_i$  le temps de séjour censuré pour l'individu  $i$  si  $\delta_i = 1$

$n_i$  le nombre de transitions pour l'individu  $i$

## ESTIMATION PARAMÉTRIQUE

- Modélisation des distributions des **temps de séjour** par des **lois paramétriques**
- Modèle à risques proportionnels pour les intensités des temps de séjour

$$\alpha_{hj}(d | \mathbf{Z}) = \alpha_{hj0}(d) \exp(\beta_{hj}^T \mathbf{Z})$$

$\alpha_{hj0}$  risque d'une loi de Weibull, Weibull généralisée, ...

$\mathbf{Z}$  vecteur des covariables

$\beta_{hj}$  vecteur des coefficients de régression

- Estimation des paramètres par maximisation de la vraisemblance



# APPLICATION À L'ASTHME

## PACKAGES DU LOGICIEL

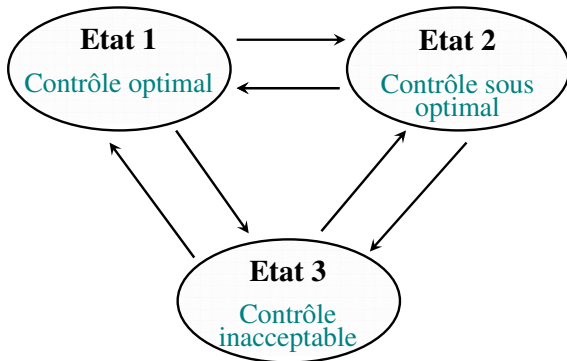
- CRAN Task View : *Survival Analysis*
- Modèle Markovien homogène : *msm*
- Modèle Markovien non homogène : *etm, mstate, msSurv*
- Modèle semi-Markovien homogène : *SemiMarkov*

## BASE DE DONNÉES

- 406 patients asthmatiques  
4.15 consultations par patient en moyenne  
Reculs de 3 mois à quatre ans
- Notion de **contrôle** : jugement global du médecin sur l'activité de la maladie
- Etude de l'**Indice de Masse Corporelle (IMC)**
  - $IMC = poids / (taille)^2$
  - 2 modalités :  $IMC < 25$  et  $IMC \geq 25$  (surpoids)

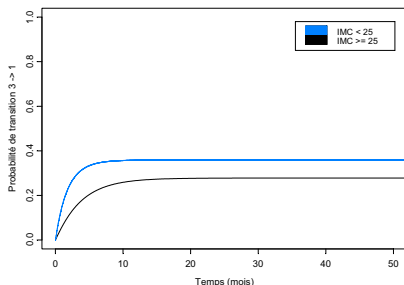
## MODÈLE

- Modèle à 3 états de contrôle
- Toutes les transitions entre les états sont possibles



## Modèle homogène

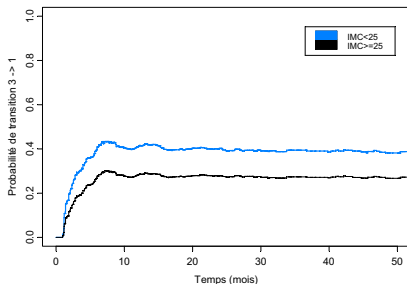
- Probabilités de transition 3  $\rightarrow$  1 pour la covariable IMC



- Modèle homogène par période (2 périodes)
  - ⇒ Meilleur ajustement aux données
  - ⇒ L'hypothèse d'homogénéité semble trop restrictive

## Modèle non-homogène

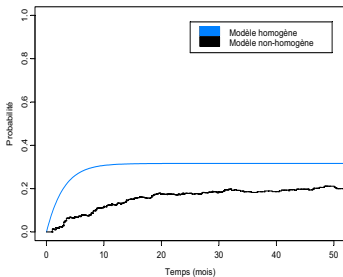
- Probabilités de transition 3  $\rightarrow$  1 pour la covariable IMC



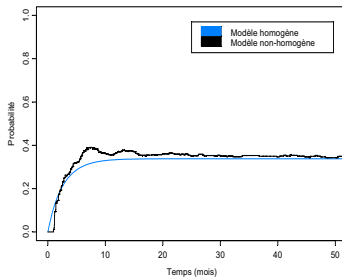
## Modèle non-homogène

- Estimations dans un modèle **homogène** et **non-homogène**

Probabilités de transition 1 vers 3

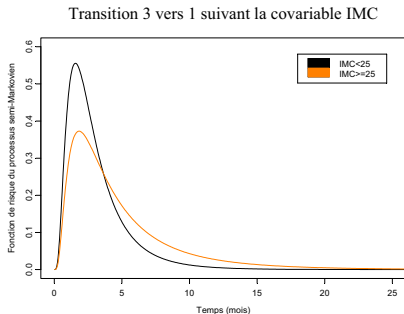


Probabilités de transition 3 vers 1



## Modèle semi-Markovien

- Estimations paramétriques des intensités  $3 \rightarrow 1$  du processus semi-Markovien pour la covariable IMC



## Coefficients de régression de l'IMC

- Coefficients de régression pour la covariable IMC

Transition	Modèle MH			Modèle MNH		
	$\beta$	<i>E.T.</i>	<i>p</i> -value	$\beta$	<i>E.T.</i>	<i>p</i> -value
1 $\rightarrow$ 2	-0.409	0.383	0.02	-0.248	0.236	0.28
1 $\rightarrow$ 3	-0.122	0.269	0.71	0.655	0.303	0.03
2 $\rightarrow$ 1	0.542	0.364	<0.01	-0.030	0.203	0.88
2 $\rightarrow$ 3	0.041	0.378	0.87	0.204	0.249	0.42
3 $\rightarrow$ 1	<b>-1.170</b>	<b>0.370</b>	<b>&lt;0.01</b>	<b>-0.595</b>	<b>0.220</b>	<b>&lt;0.01</b>
3 $\rightarrow$ 2	-0.561	0.293	<0.01	-0.234	0.191	0.21

- Effet de l'IMC également significatif pour la transition 3  $\rightarrow$  1
  - par stratification
  - avec l'IMC dépendant du temps
  - avec le modèle semi-Markovien




## Coefficients de régression de l'IMC

- Modèle Markov homogène avec 2 états de contrôle

Transition	Covariables	Modèle univarié			Modèle multivarié		
		$\hat{\beta}$	<i>E.T.</i>	<i>p</i> -value	$\hat{\beta}$	<i>E.T.</i>	<i>p</i> -value
1 → 2	IMC	-0.129	0.247	0.60	-0.174	0.278	0.53
	Sévérité	0.665	0.272	0.04	0.820	0.305	<0.01
	Corticoïdes Oraux	0.110	0.269	0.68	-0.422	0.305	0.17
	Antécédents Corticoïdes	0.651	0.262	0.01	0.498	0.299	0.10
2 → 1	IMC	<b>-0.801</b>	<b>0.184</b>	<b>&lt;0.01</b>	<b>-0.637</b>	<b>0.219</b>	<b>&lt;0.01</b>
	Sévérité	-0.726	0.203	<0.01	-0.062	0.255	0.81
	Corticoïdes Oraux	-1.002	0.209	<0.01	-0.693	0.248	<0.01
	Antécédents Corticoïdes	-0.852	0.212	<0.01	-0.312	0.266	0.24

⇒ Résultats ajustés

# DISCUSSION

- Méthodes adaptées pour modéliser des données de survie multivariée
- Résultats facilement interprétables
- Packages 
- Hypothèses sous-jacentes (Markov, homogénéité, risques proportionnels, observation continue)
- Dépendances entre individus : effets aléatoires