

# Analyse de données longitudinales : Application aux jeux vidéo

Young Statisticians and Probabilists

Thibault ALLART

Encadrants  
de thèse

Agathe  
Guilloux

Stéphane  
Natkin

Guillaume  
Levieux

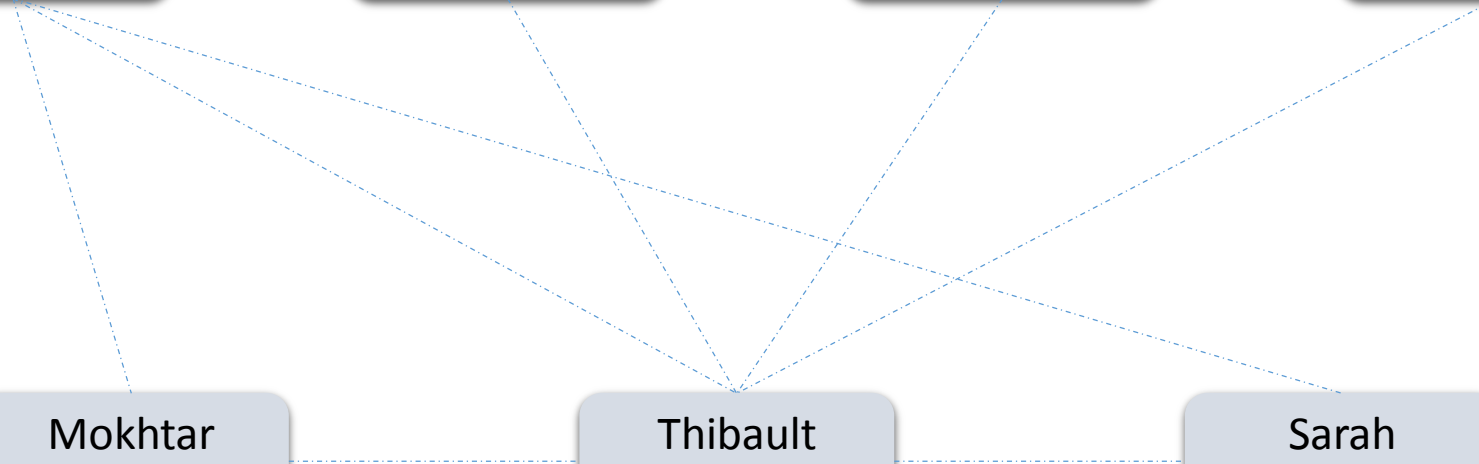
Michel  
Pierfitte

Collaborations

Mokhtar  
Zahdi Alaya

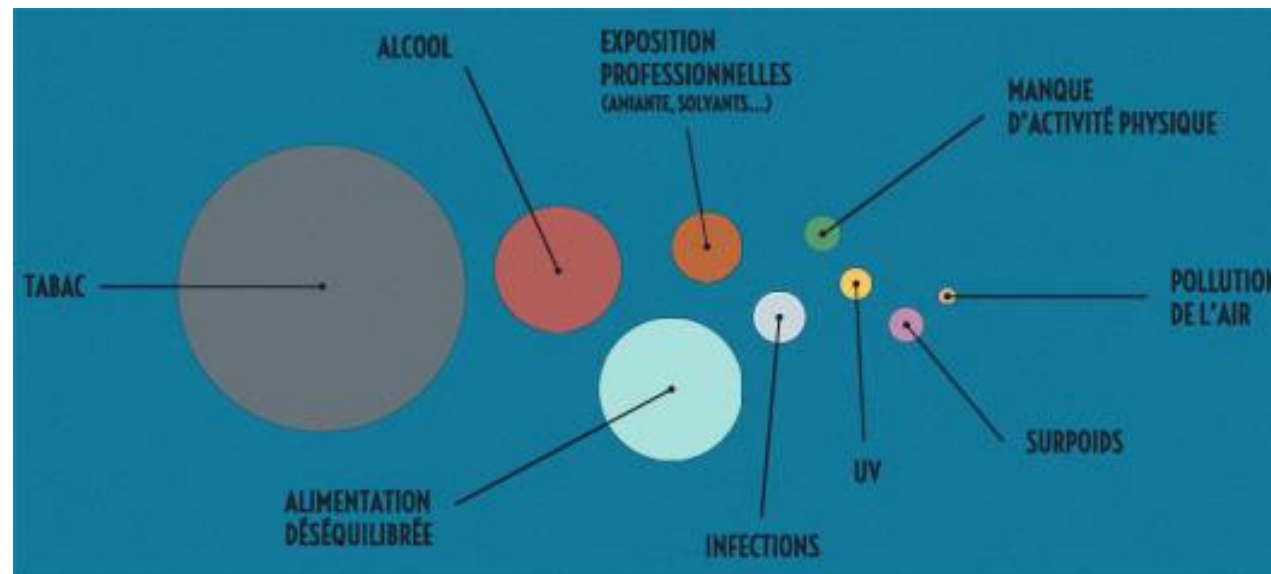
Thibault  
ALLART

Sarah  
Lemler



On s'intéresse à l'influence de certains facteurs sur le temps d'apparition d'un évènement.

## Poids des différents facteurs de risque de cancer



Source : institut national du cancer

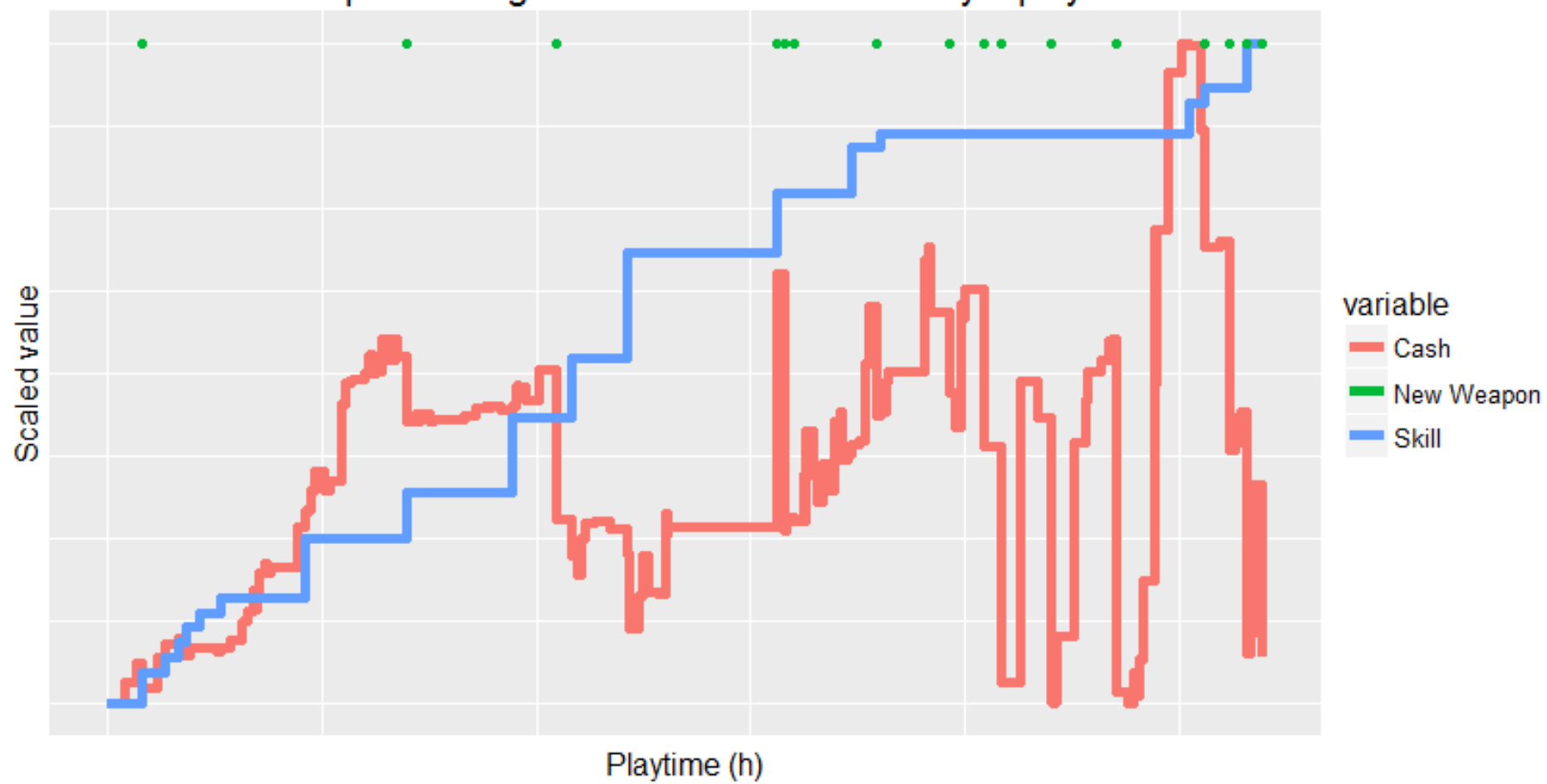
Notre exposition change  
au cours du temps

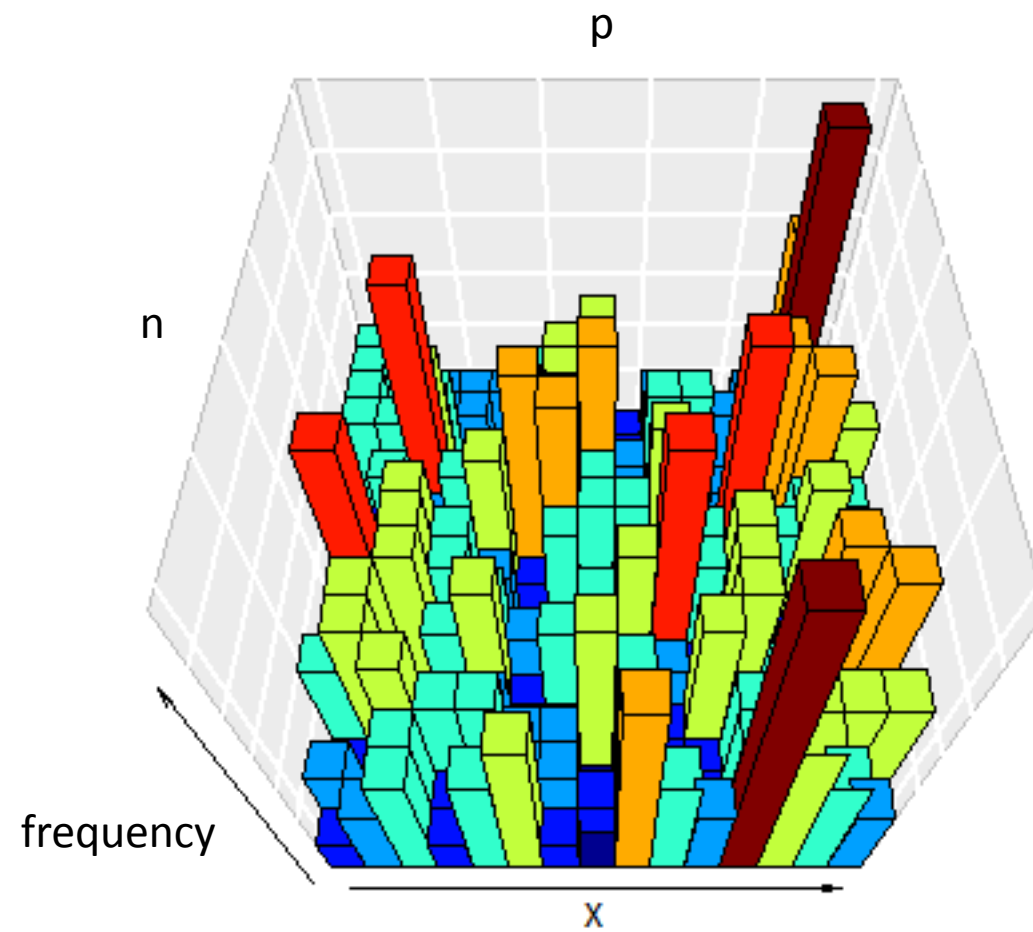


→  
Playtime

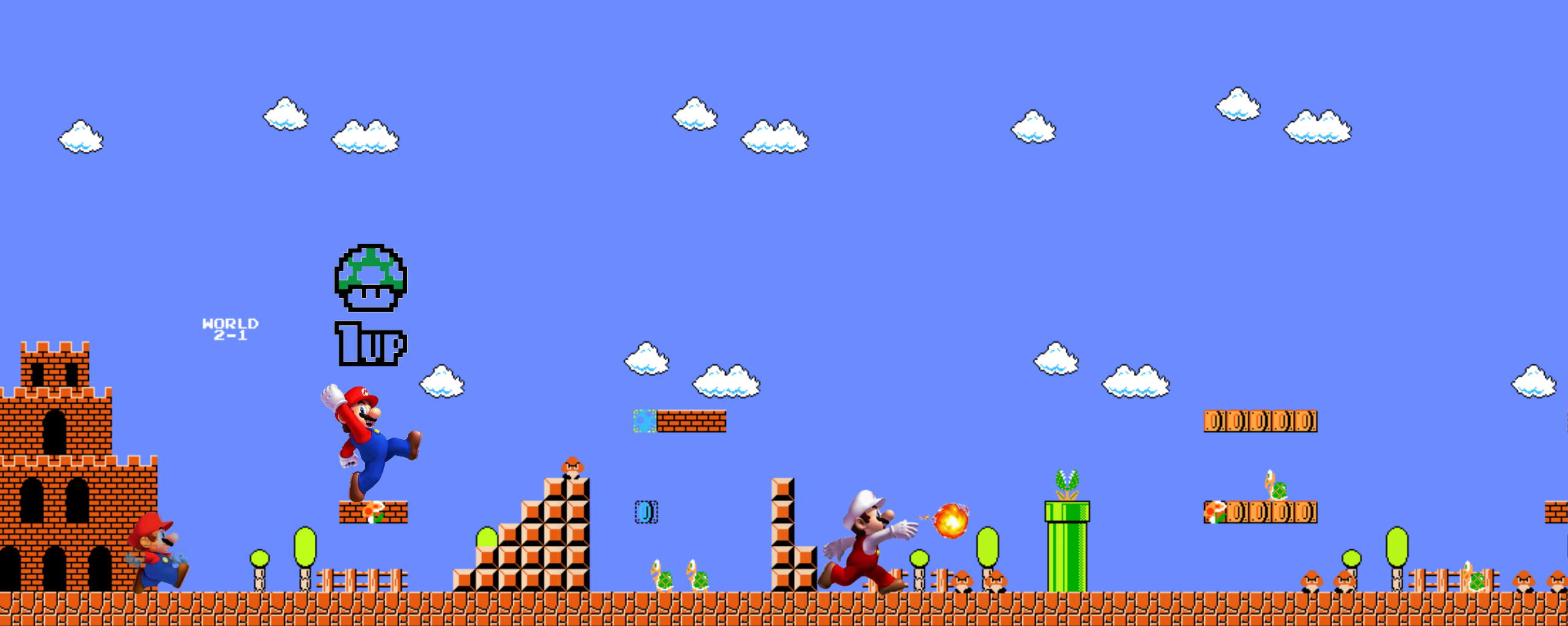
Longitudinal data

Example of Longitudinal Data for one Far Cry 4 player





Observation frequency depends on individuals **and** covariates



→  
Playtime

Why do we need Longitudinal data ?



Player 1



Quits before  
the end

	Retention
Fire Mario	Correlated
Water	Correlated

P1 Stop

2h

P2 Stop

30h

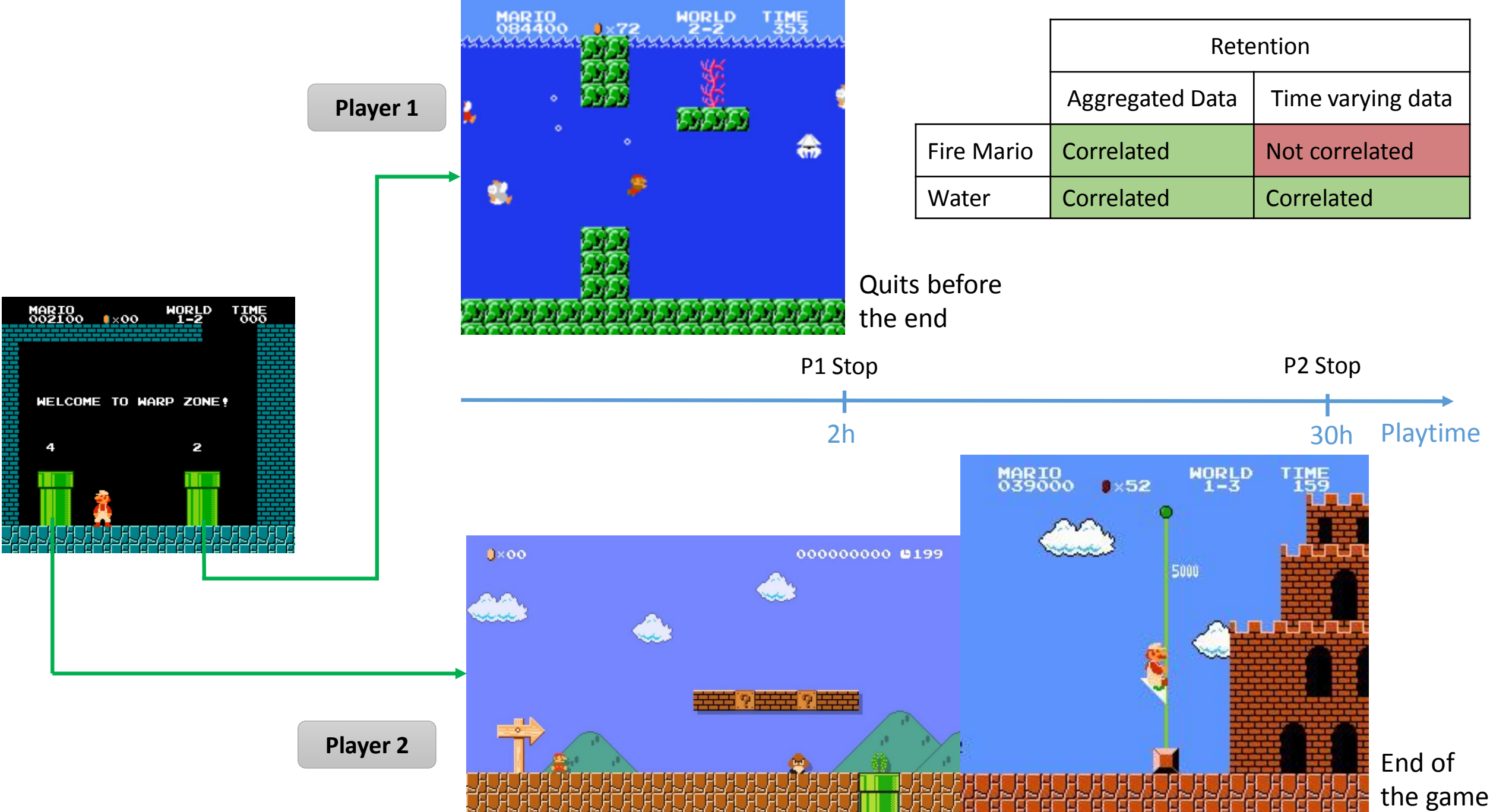
Playtime

Player 2



End of  
the game





Idem : on ne peut pas comparer les indicateurs de santé d'un nouveau né avec ceux d'un adulte

Joueur 1



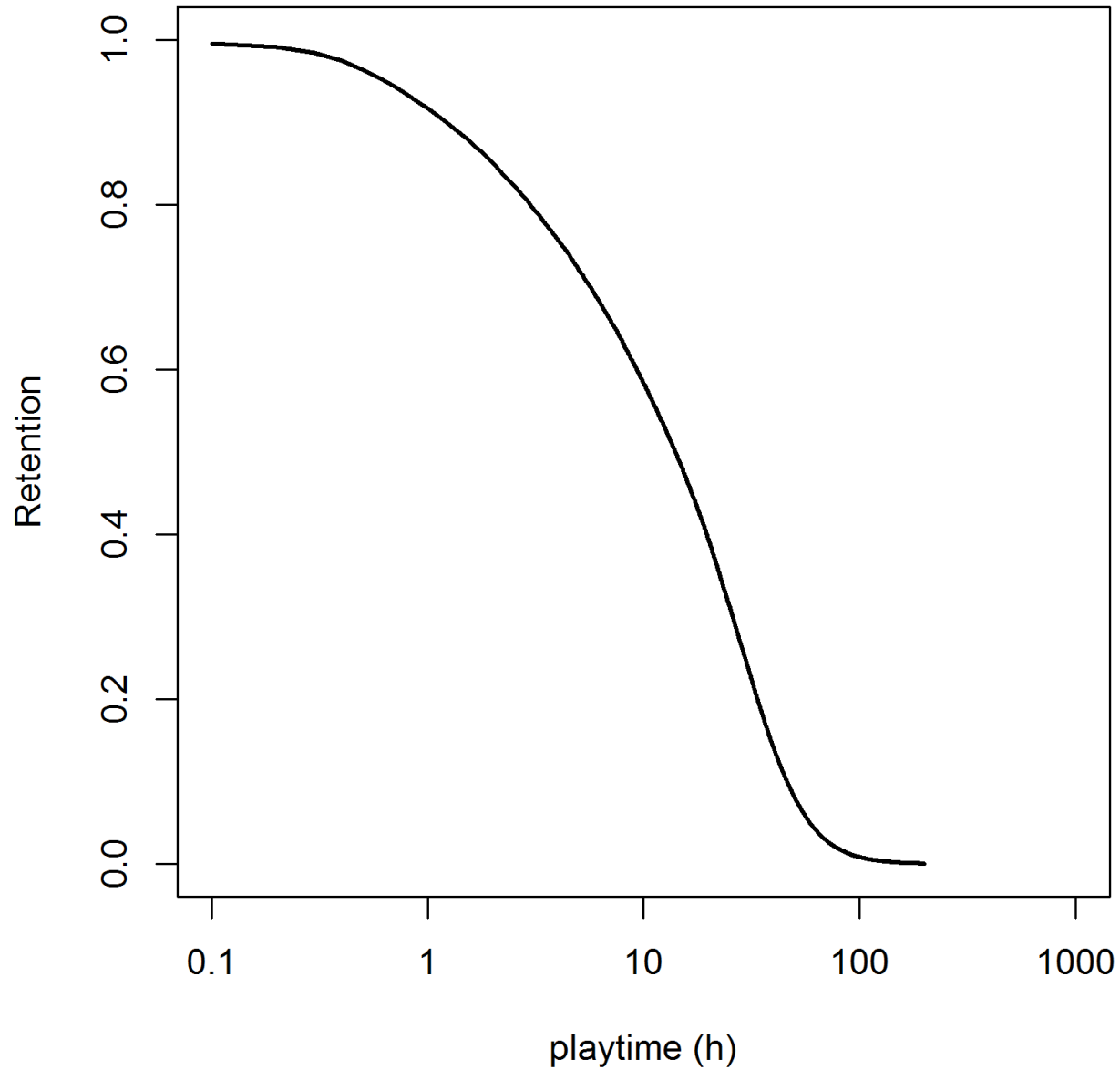
Joueur 2



# Modélisation

Analyse de survie

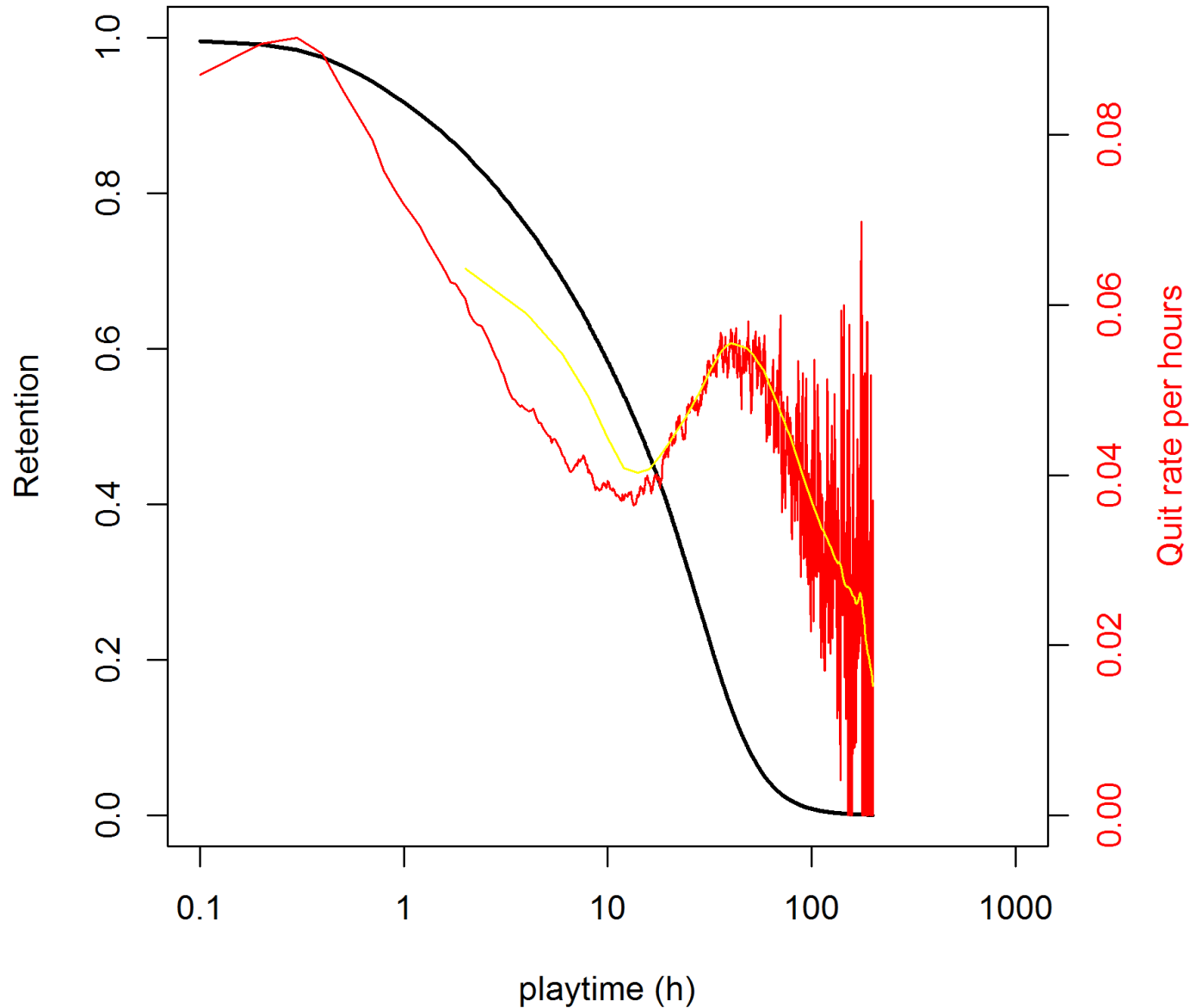
## Far Cry 4 retention



Soit  $T$  la variable aléatoire positive associée au temps de jeu (playtime) des joueurs.

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

## Far Cry 4 quit rate



Soit  $T$  la variable aléatoire positive associée au temps de jeu (playtime) des joueurs.

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

Taux de risqué instantané :

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{\mathbb{P}(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

On cherche à modéliser  $\lambda(t)$

# Time-varying Cox proportional hazard model

$$\lambda(t|X_i(t)) = \lambda_0(t) \exp \left( \sum_{j=1}^p X_{i,j}(t) \beta_j(t) \right)$$

Diagram illustrating the components of the Time-varying Cox proportional hazard model equation:

- Quit rate** (red text) points to  $\lambda(t|X_i(t))$ .
- Baseline** (purple text) points to  $\lambda_0(t)$ .
- Player behavior** (green text) points to  $X_{i,j}(t)$ .
- Effect of player behavior** (orange text) points to  $\beta_j(t)$ .

Comment estimer les coefficients du modèle ?

# Estimation

Maximum de vraisemblance

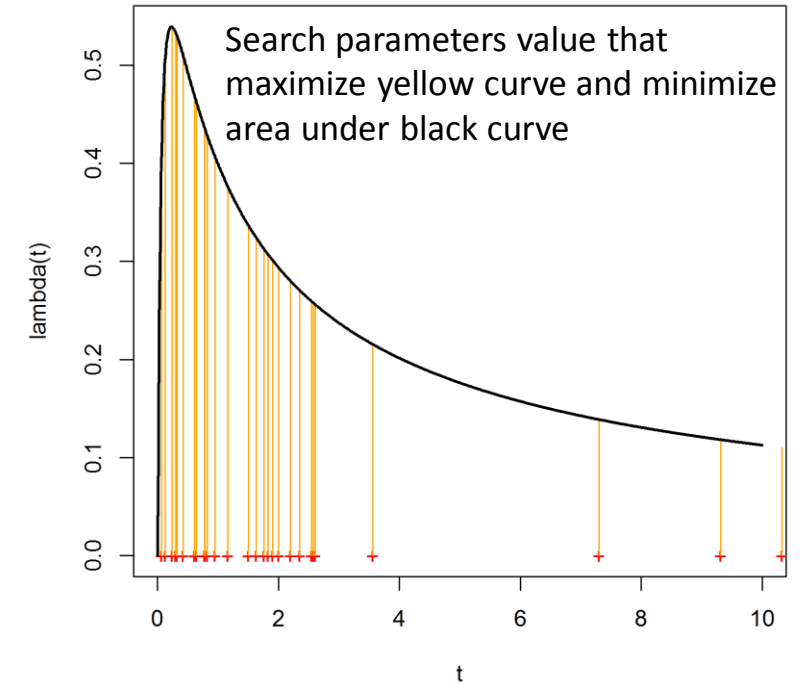
# Likelihood

Soit  $t_i$  le temps de réalisation de l'évènement pour l'individu  $i$  et  $\tau$  le temps de fin d'observation, alors la vraisemblance se décompose comme suit :

Probabilité que les évènements aient eu lieu au moment où on les a observés, conditionnellement au passé du processus

$$\mathcal{L}_n(\beta) = \left\{ \prod_{t_i \leq \tau} \lambda(t_i) \right\} S(\tau)$$

Probabilité qu'il n'y ait pas d'évènement aux autres temps, conditionnellement au passé du processus





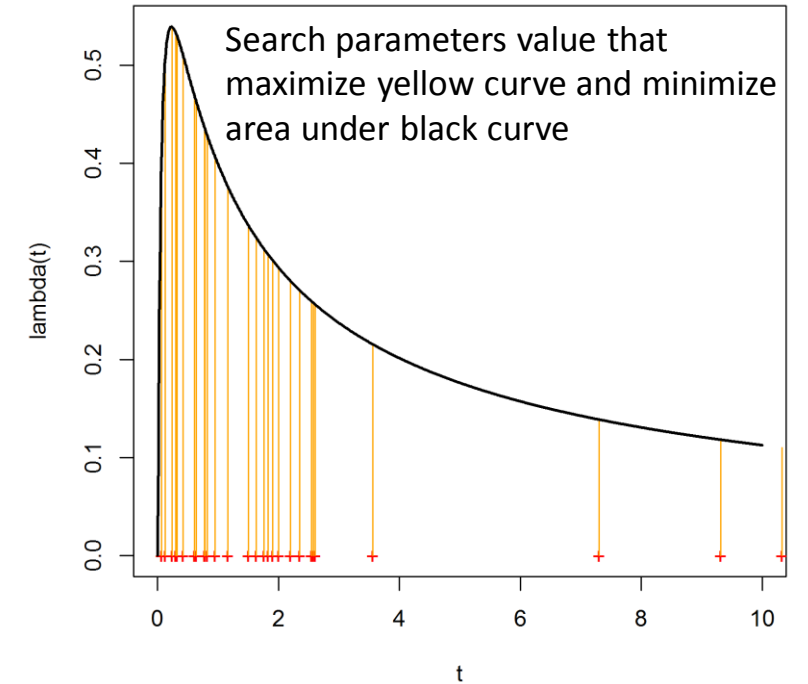
# Likelihood

Soit  $\tau_i$  le temps de réalisation de l'évènement pour l'individu  $i$ , alors la vraisemblance se décompose comme suit :

Probabilité que les évènements aient eu lieu au moment où on les a observés, conditionnellement au passé du processus

Probabilité qu'il n'y ait pas d'évènement aux autres temps, conditionnellement au passé du processus

$$\mathcal{L}_n(\beta) = \left\{ \prod_{t_i \leq \tau} \lambda(t_i) \right\} \exp \left( - \int_0^{\tau} \lambda(s) ds \right)$$



# Likelihood

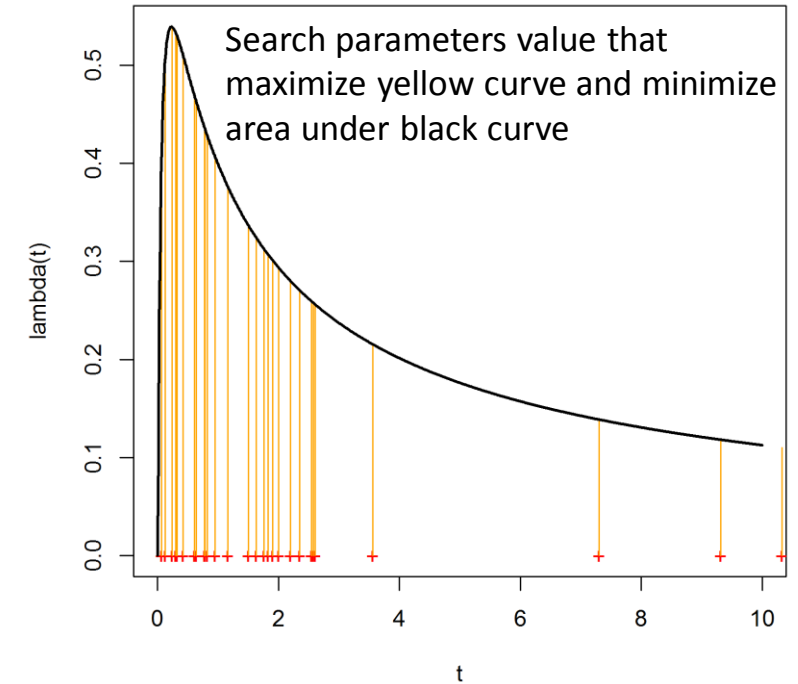
Soit  $\tau_i$  le temps de réalisation de l'évènement pour l'individu  $i$ , alors la vraisemblance se décompose comme suit :

Probabilité que les évènements aient eu lieu au moment où on les a observés, conditionnellement au passé du processus

Probabilité qu'il n'y ait pas d'évènement aux autres temps, conditionnellement au passé du processus

$$\mathcal{L}_n(\beta) = \left\{ \prod_{t_i \leq \tau} \lambda(t_i) \right\} \exp \left( - \int_0^{\tau} \lambda(s) ds \right)$$

$$\mathcal{L}_n(\beta) = \left\{ \prod_{t_i \leq \tau} \exp(X^T(t_i)\beta(t_i)) \right\} \cdot \exp \left( - \int_0^{\tau} Y(s) \exp(X^T(s)\beta(s)) ds \right)$$



# Likelihood

Soit  $\tau_i$  le temps de réalisation de l'évènement pour l'individu  $i$ , alors la vraisemblance se décompose comme suit :

Probabilité que les évènements aient eu lieu au moment où on les a observés, conditionnellement au passé du processus

Probabilité qu'il n'y ait pas d'évènement aux autres temps, conditionnellement au passé du processus

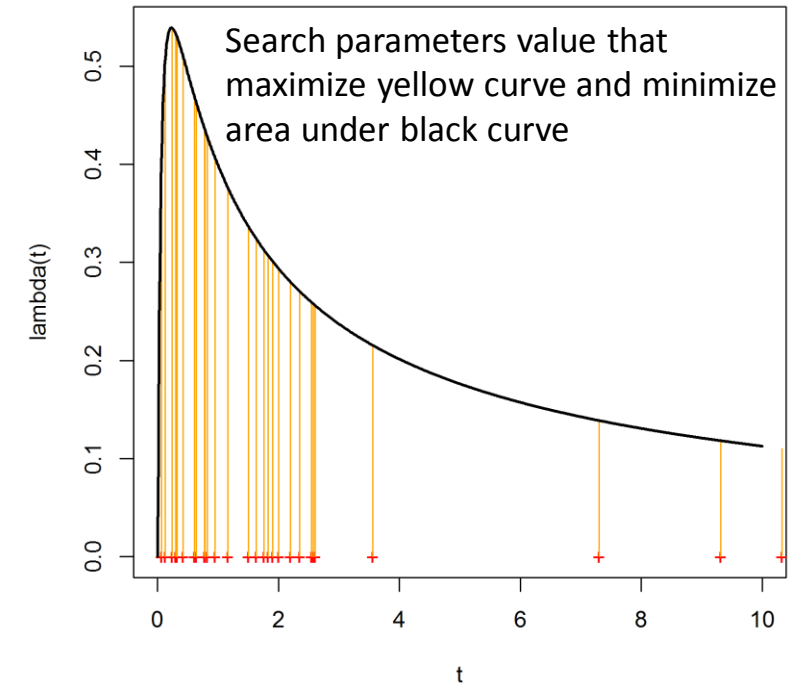
$$\mathcal{L}_n(\beta) = \left\{ \prod_{t_i \leq \tau} \lambda(t_i) \right\} \exp \left( - \int_0^{\tau} \lambda(s) ds \right)$$

$$\mathcal{L}_n(\beta) = \left\{ \prod_{t_i \leq \tau} \exp \left( X^T(t_i) \beta(t_i) \right) \right\} \cdot \exp \left( - \int_0^{\tau} Y(s) \exp \left( X^T(s) \beta(s) \right) ds \right)$$

La log vraisemblance s'écrit :

$$\ell_n(\beta) = \left\{ \sum_{t_i \leq \tau} X^T(t_i) \beta(t_i) \right\} - \int_0^{\tau} Y(s) \exp \left( X^T(s) \beta(s) \right) ds$$

$$\ell_n(\beta) = \frac{1}{n} \sum_{i=1}^n \left\{ \int_0^{\tau} X_i^T(t) \beta(t) dN_i(t) - \int_0^{\tau} Y_i(t) \exp \left( X_i^T(t) \beta(t) \right) dt \right\}$$



- Integrale
- Quelle forme pour les  $\beta(t)$  ?

# Coefficients constants par morceaux

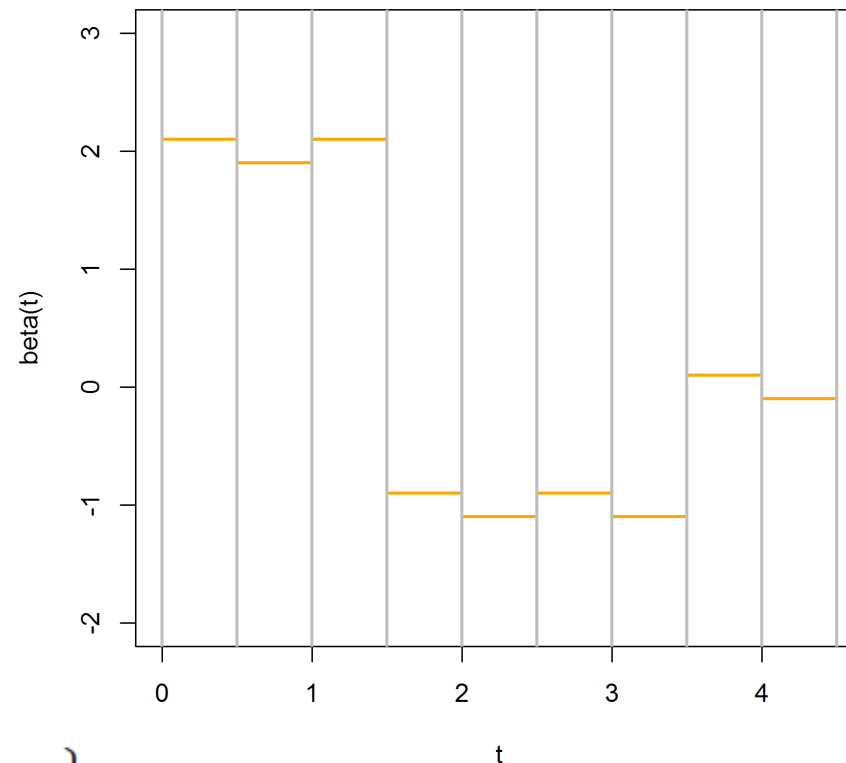
Soit  $(I_l)_{l \in \{0, L\}}$  une partition de  $[0, \tau]$ .

$$\beta_j(t) = \sum_{l=1}^L \beta_{j,l} \mathbf{1}_{(I_l)}(t)$$

Et en utilisant le fait que les  $X_{i,j}(t)$  sont constants sur de petits intervalles.

$$\ell_n(\beta) = \frac{1}{n} \sum_{i=1}^n \left\{ \int_0^\tau X_i^T(t) \beta(t) dN_i(t) - \int_0^\tau Y_i(t) \exp(X_i^T(t) \beta(t)) dt \right\}$$

$$\ell_n(\beta) = \frac{1}{n} \sum_{i=1}^n \sum_{l=1}^L \left\{ \sum_{j=1}^p \sum_s X_{i,s}^j \beta_{j,l} N_i(I_s) - \sum_s Y_i(I_s) \exp \left( \sum_{j=1}^p X_{i,s}^j \beta_{j,l} \right) |I_s| \right\}$$



# Pénalité

Pour éviter l'explosion du nombre de dimensions, on voudrait pénaliser le nombre de coefficients (ou le nombre de sauts).

$$\sum_{j=2}^p (\beta_{j+1} \neq \beta_j)$$

Cette norme n'étant pas convexe, on utilise sa relaxation convexe (Total variation ou Fused Lasso).

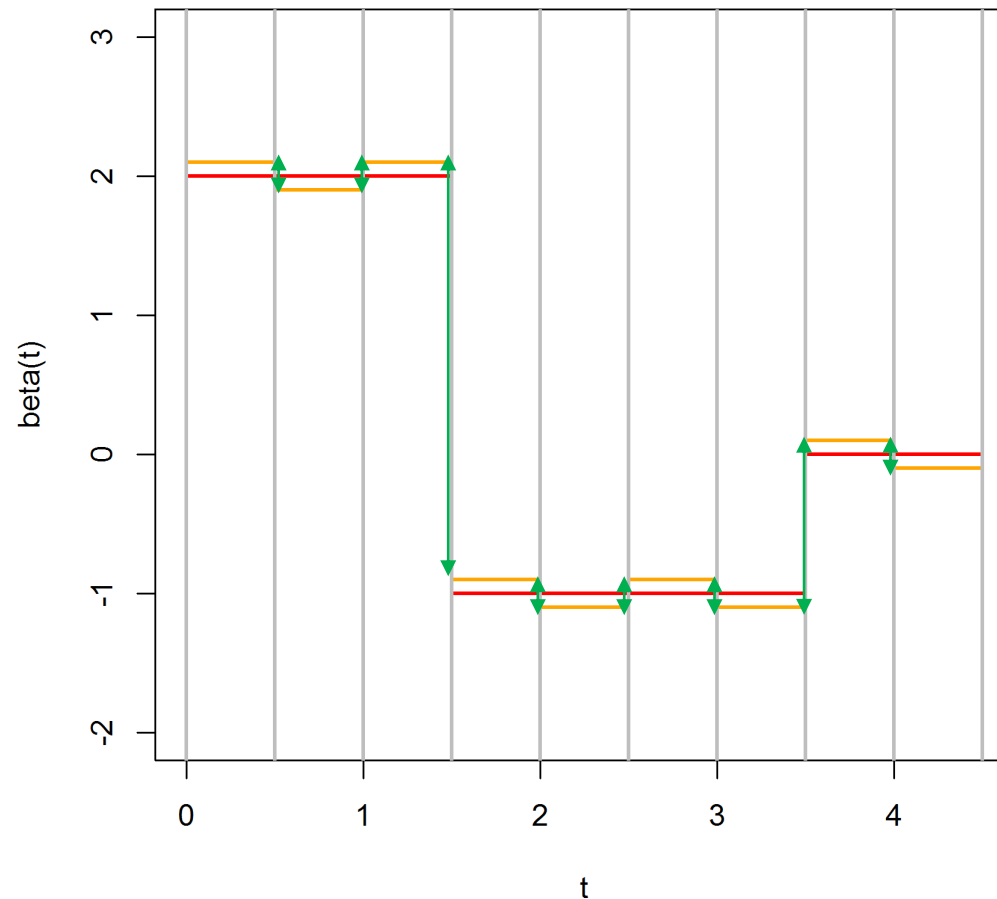
On pénalise la hauteur de sauts.

$$\|\beta\|_{\text{TV}} = \sum_{l=2}^L |\beta_{j,l} - \beta_{j,l-1}|$$

Version multivariée (group-TV):

$$\|\beta\|_{\text{TV}} = \sum_{j=1}^p \left( |\beta_{j,1}| + \sum_{l=2}^L |\beta_{j,l} - \beta_{j,l-1}| \right)$$

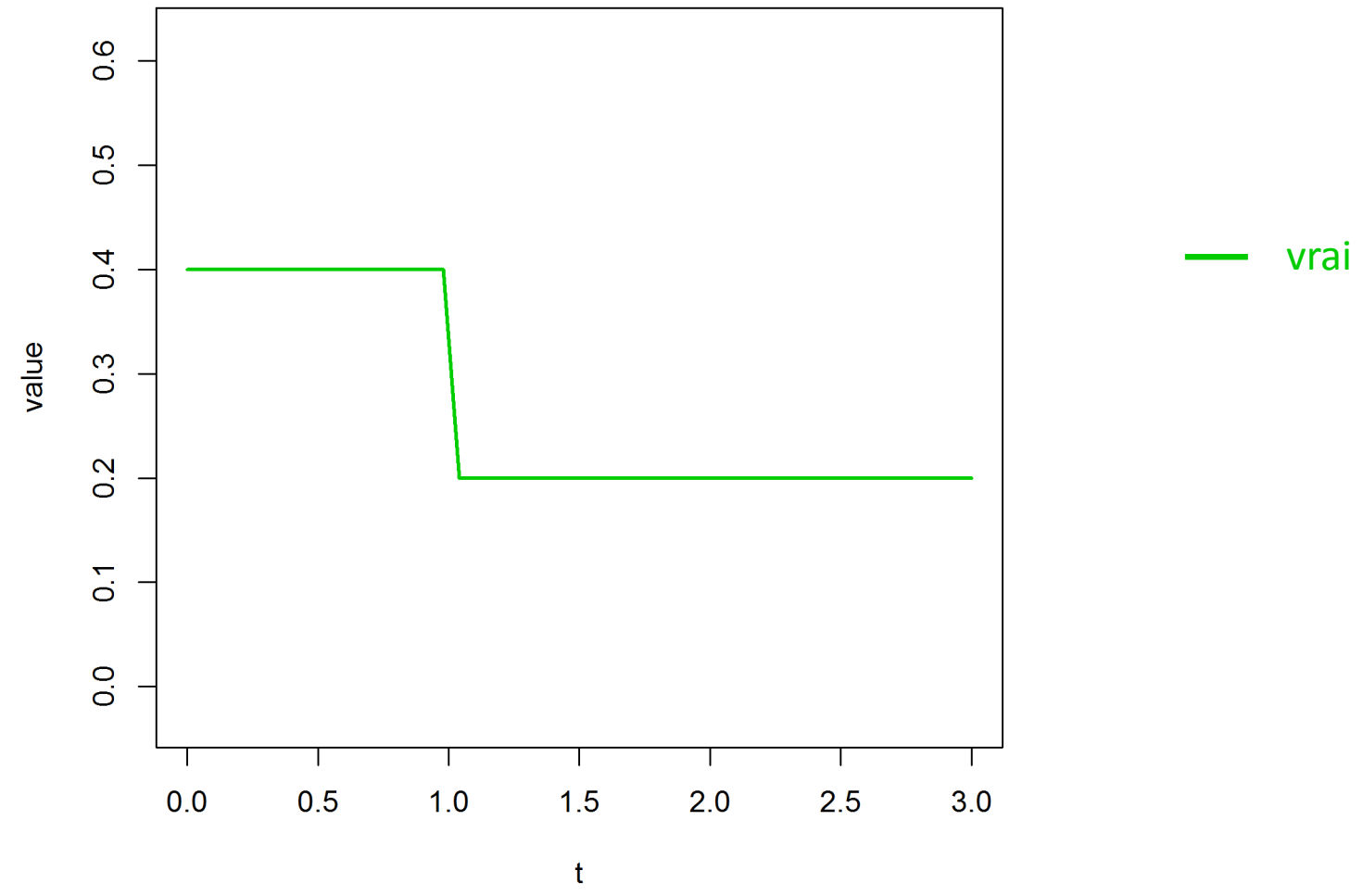
$$\hat{\beta}^{\text{TV}} \in \operatorname{argmin}_{\beta \in \mathbb{R}^L} \left\{ \hat{R}(\beta) + \lambda \|\beta\|_{\text{TV}} \right\}, \lambda > 0$$



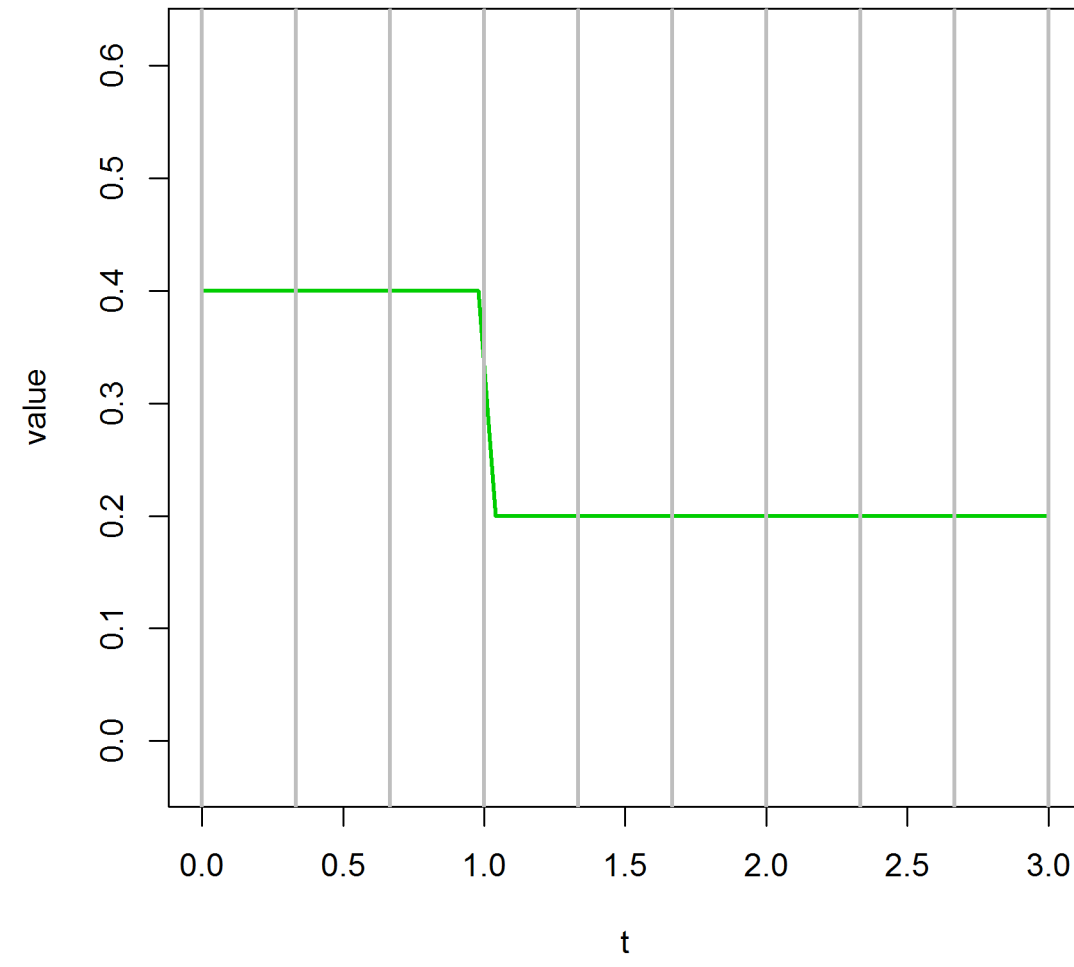
# Simulation

Thinning

**Beta 1**



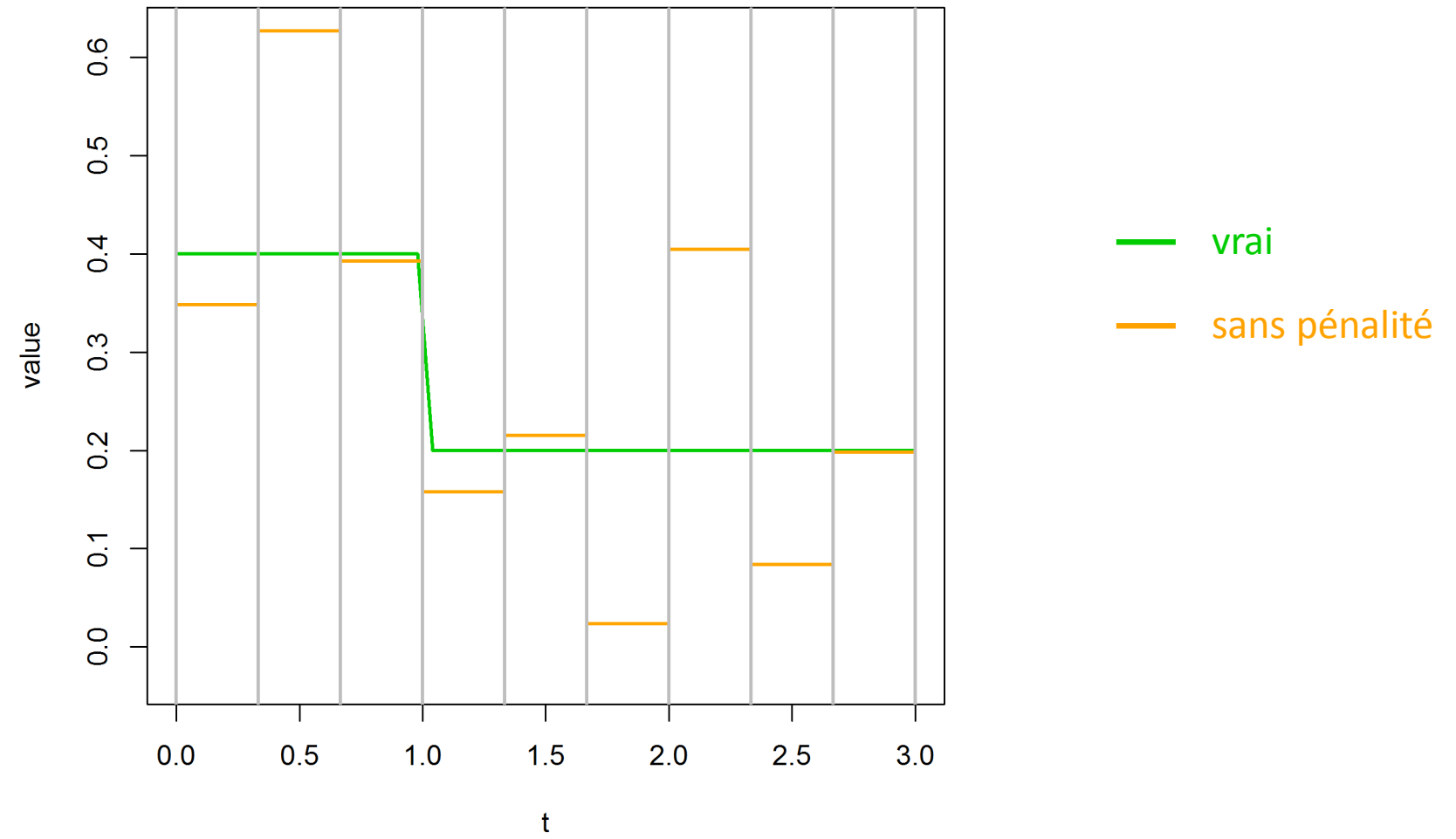
**Beta 1**



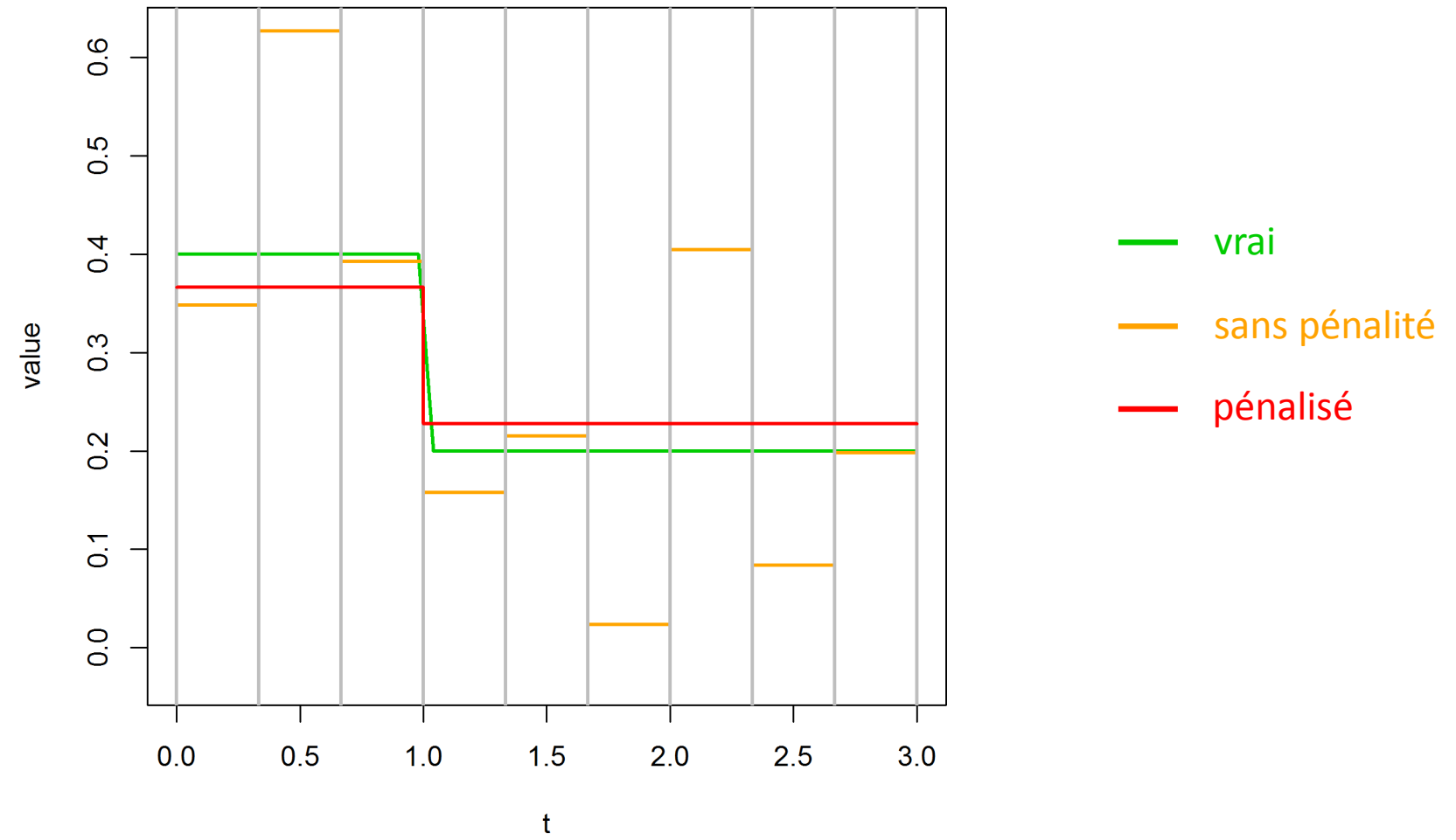
— vrai



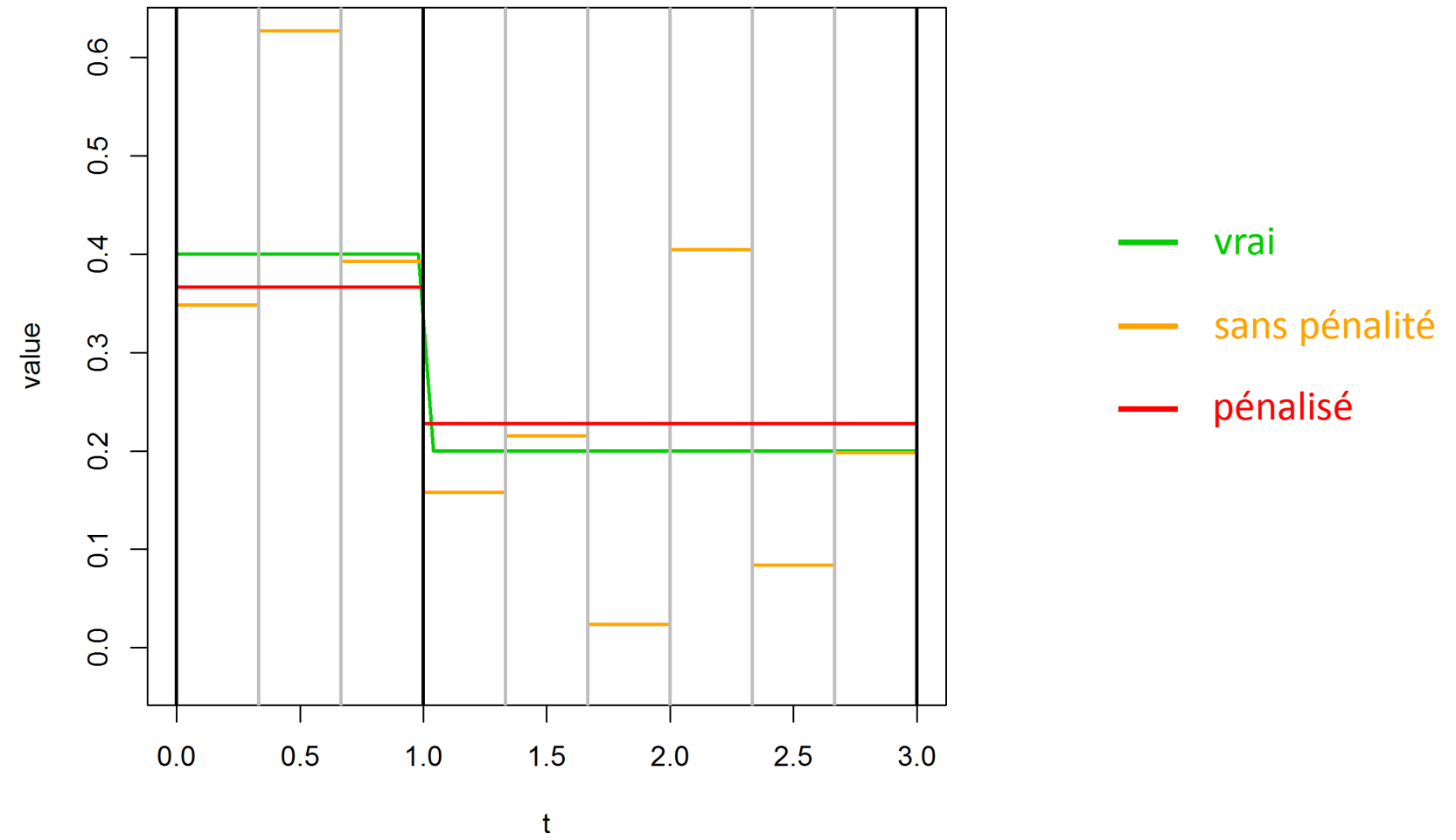
Beta 1



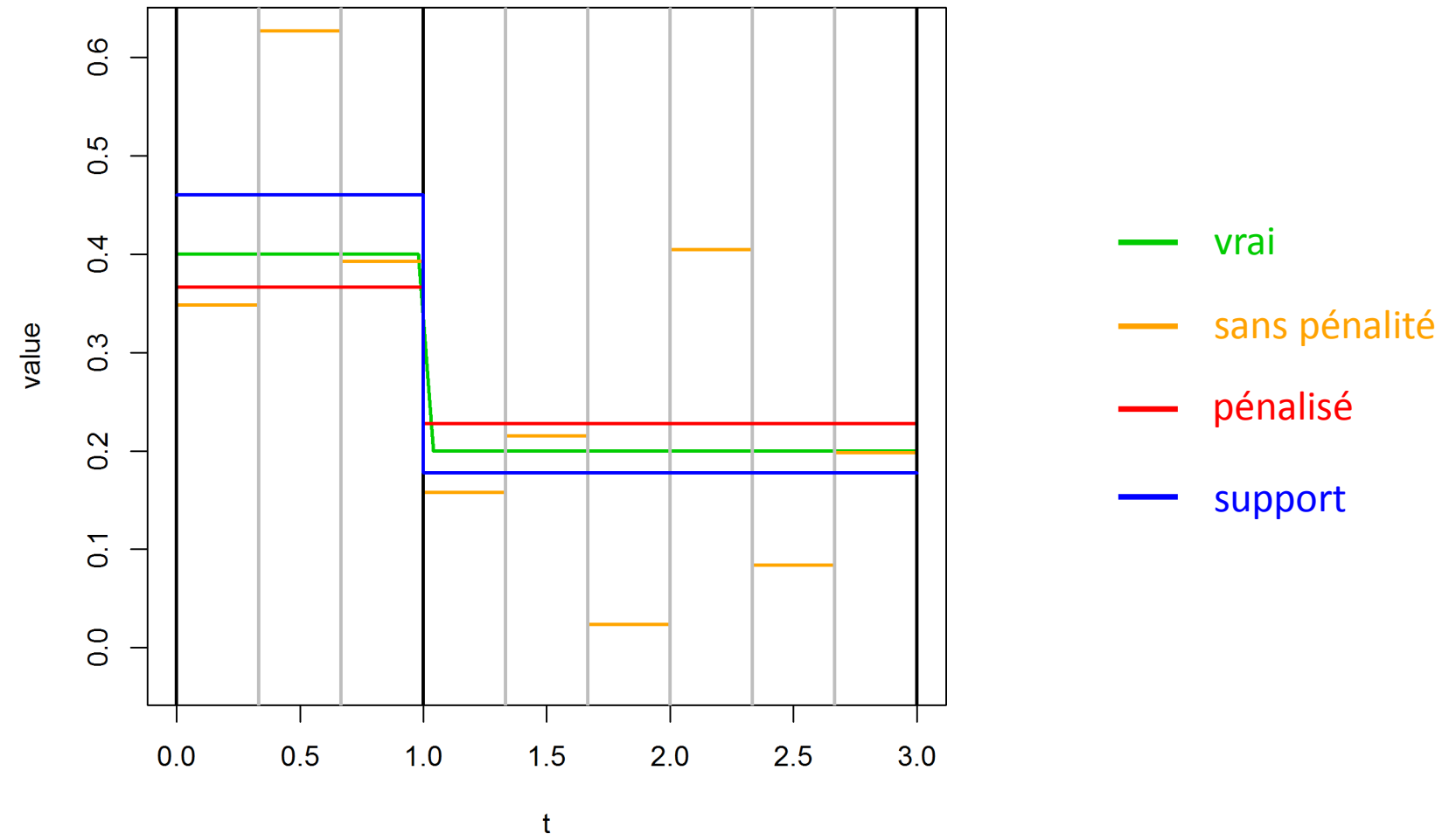
Beta 1



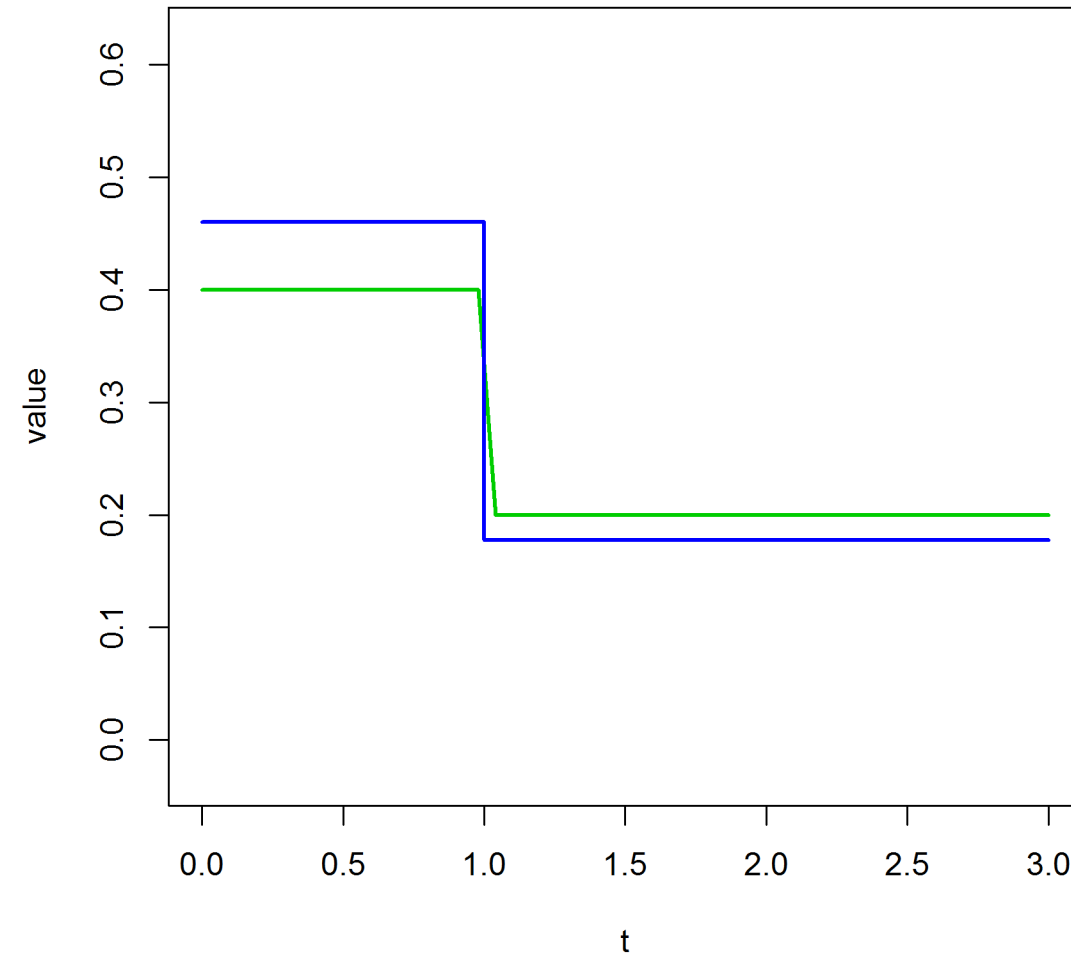
Beta 1



Beta 1



**Beta 1**



# timereg

Il existe déjà un package R pour les modèles de survie dépendant du temps.

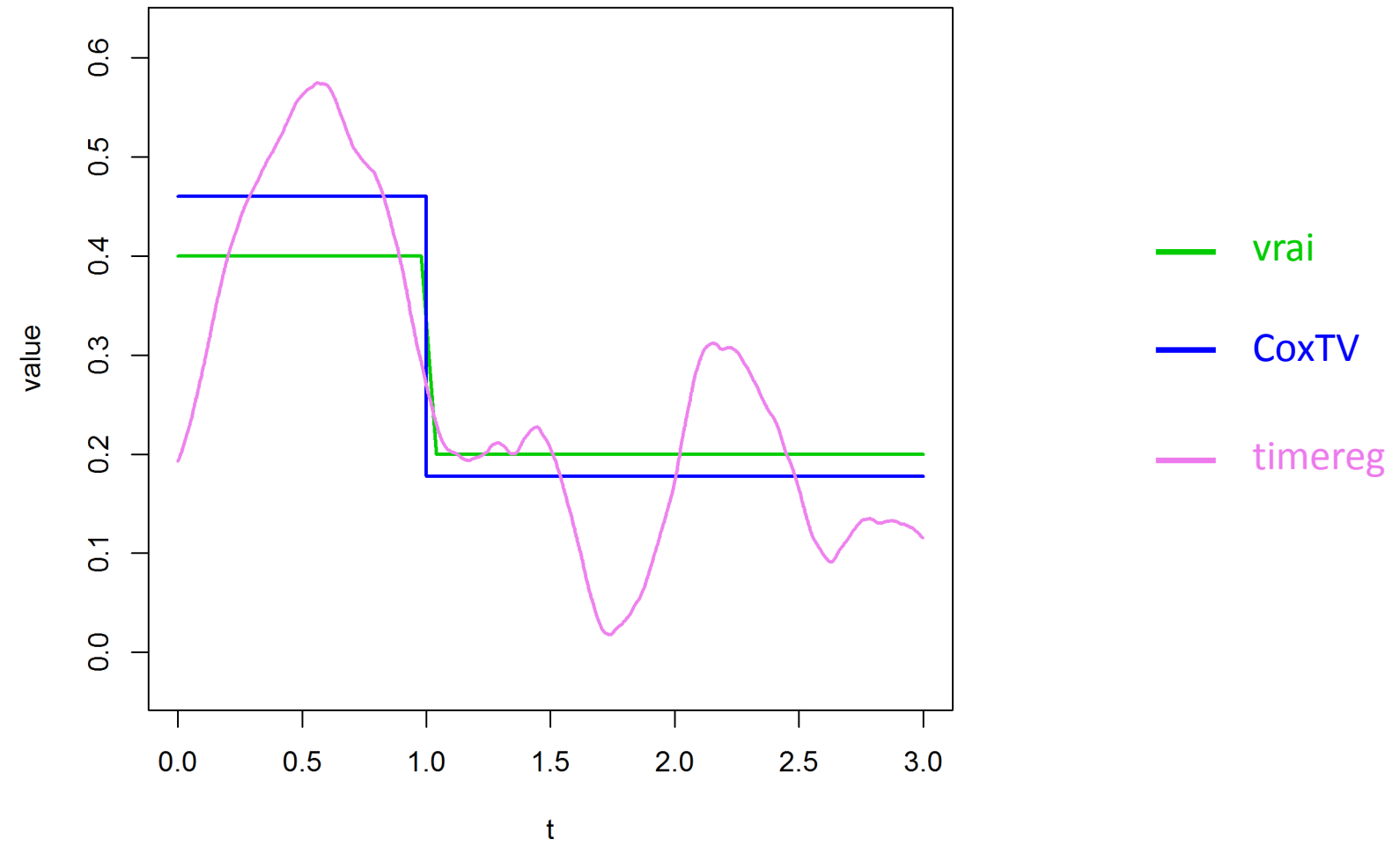
**timereg: Flexible Regression Models for Survival Data**

Programs for Martinussen and Scheike (2006), 'Dynamic Regression Models for Survival Data'

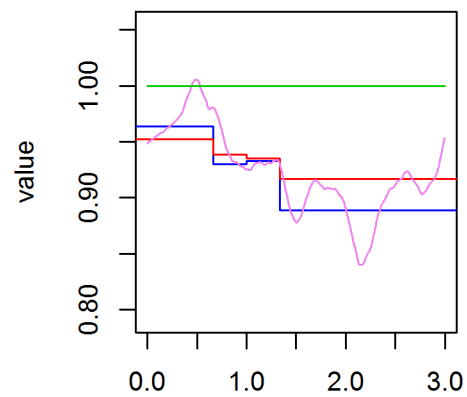
Différences :

- Estimation basée sur la fonction de hasard **cumulée**
- Retourne les coefficients cumulés, il faut ajouter un estimateur à noyaux pour obtenir les coefficients.
- L'optimisation repose sur une inversion matricielle et des itérés de lissage par noyaux

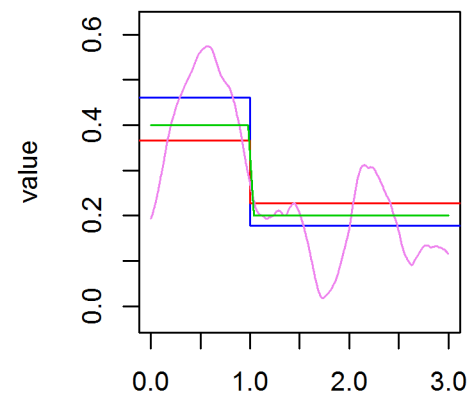
**Beta 1**



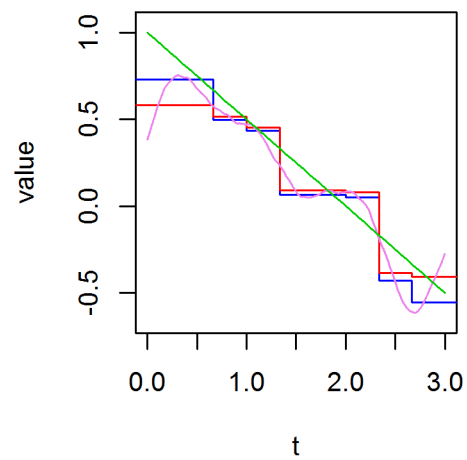
**Baseline**



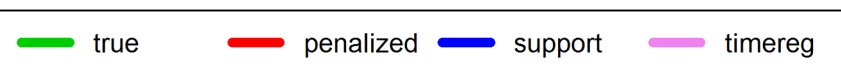
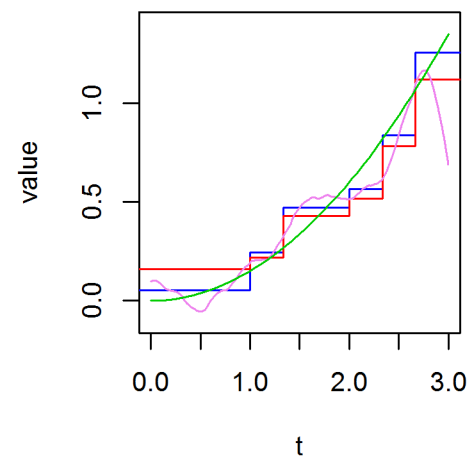
**Beta 1**



**Beta 2**

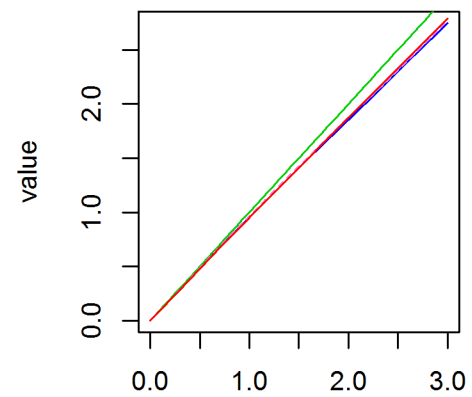


**Beta 3**

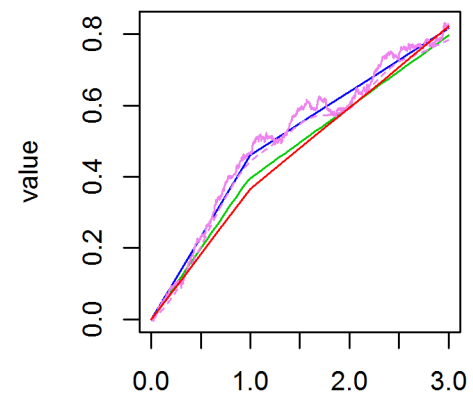




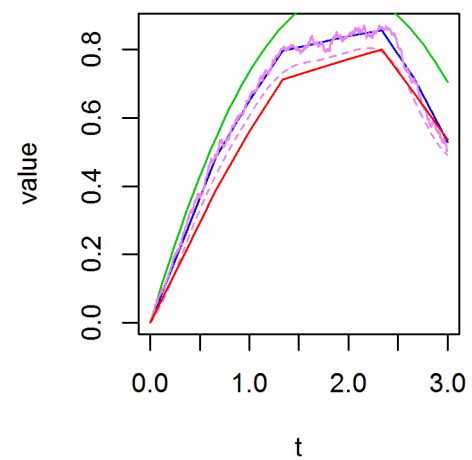
**Baseline**



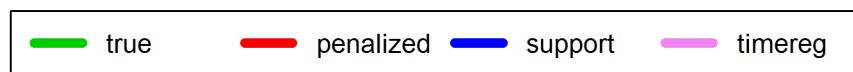
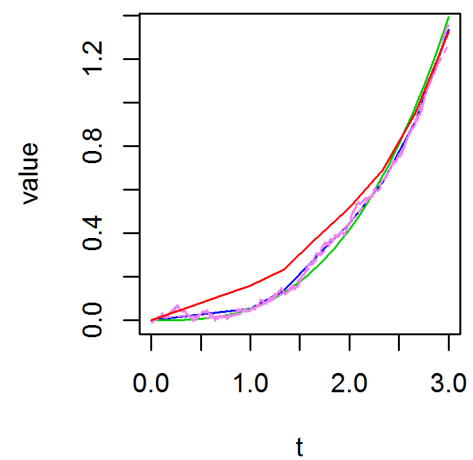
**Beta 1**



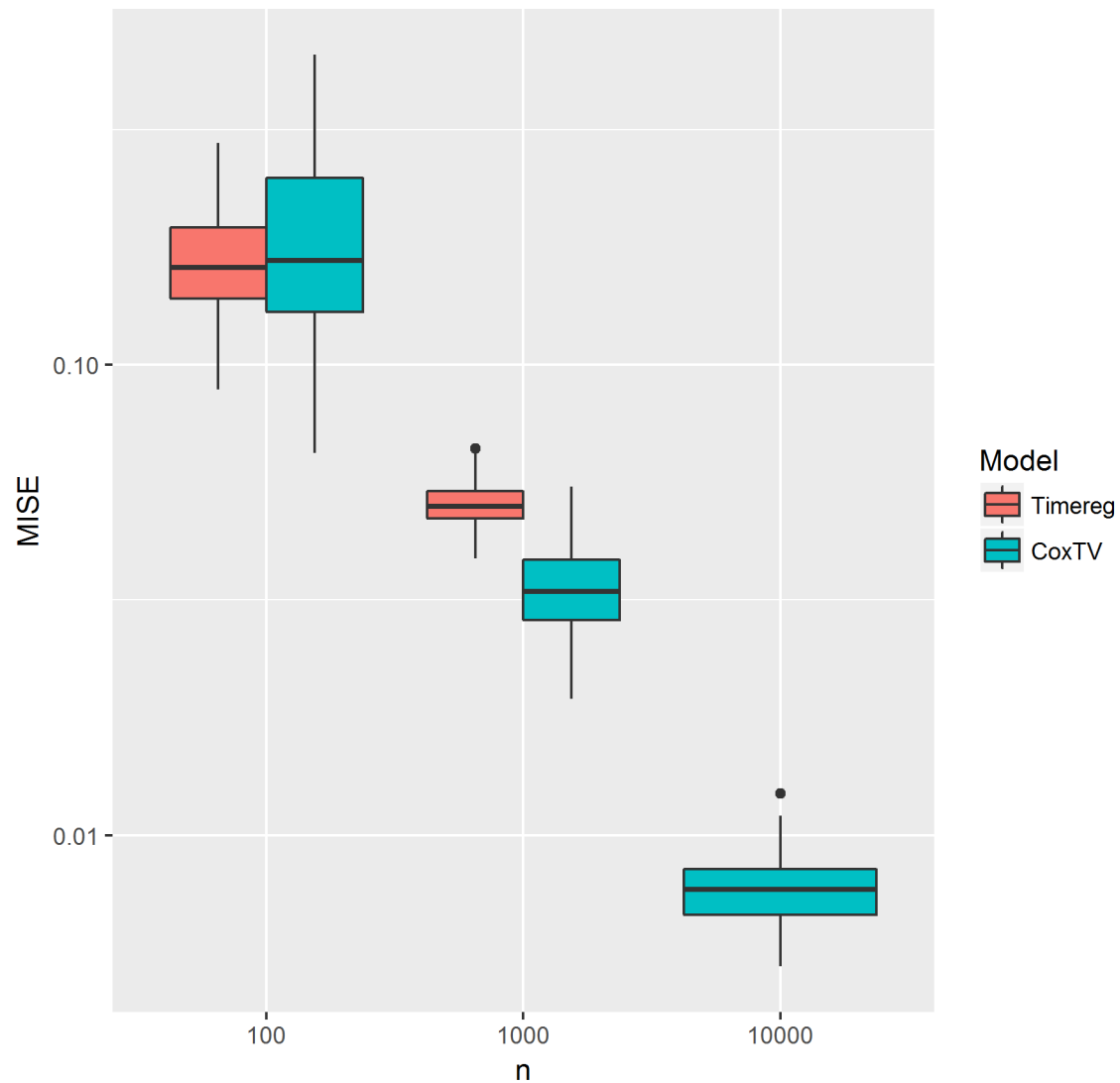
**Beta 2**



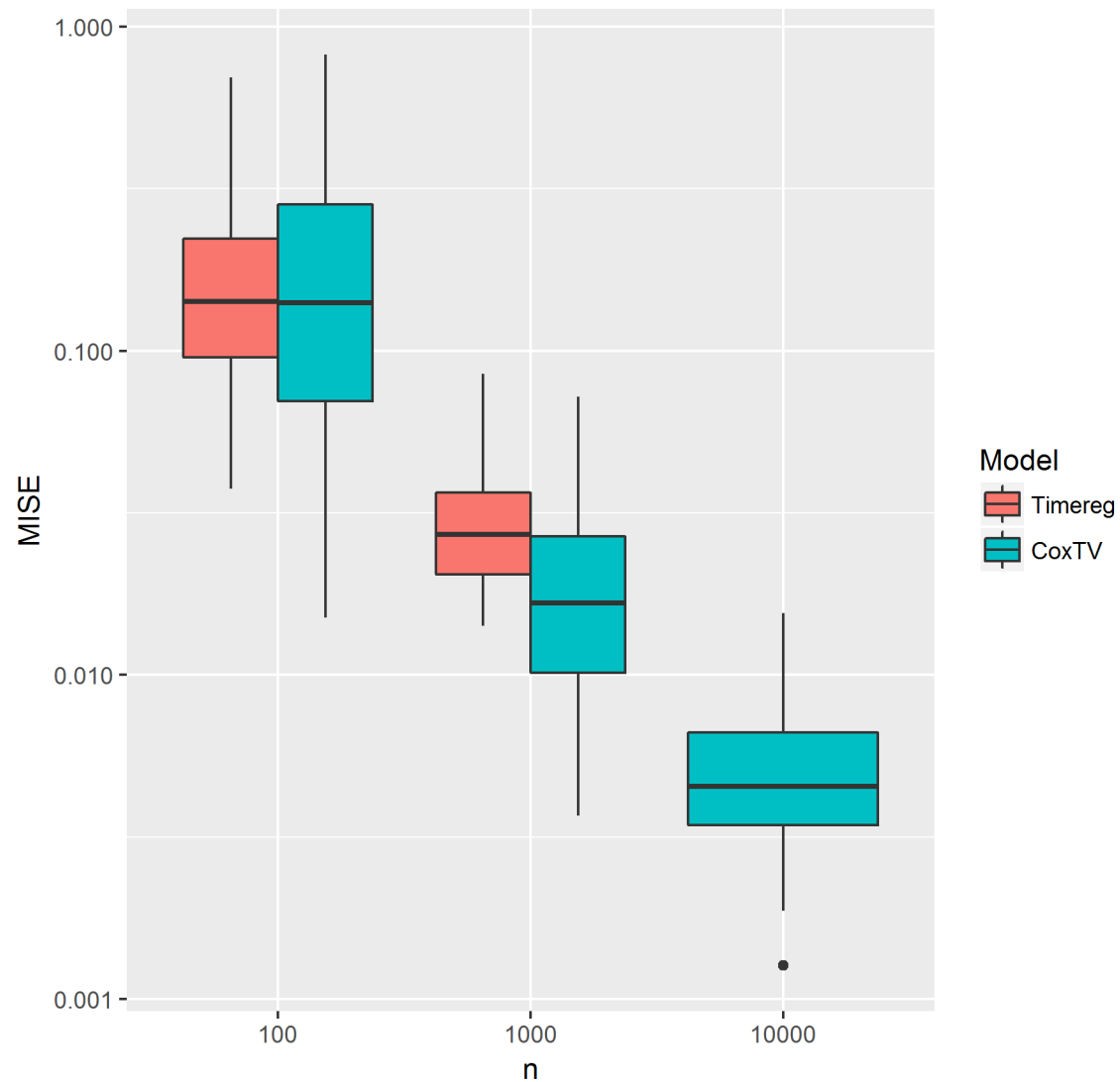
**Beta 3**



MISE on coefficients



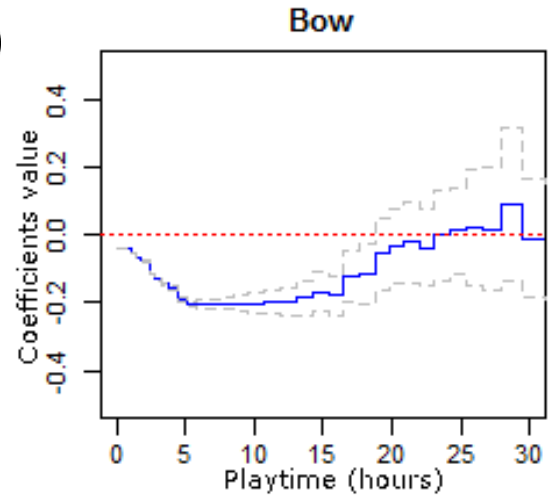
MISE on cumulated coefficients



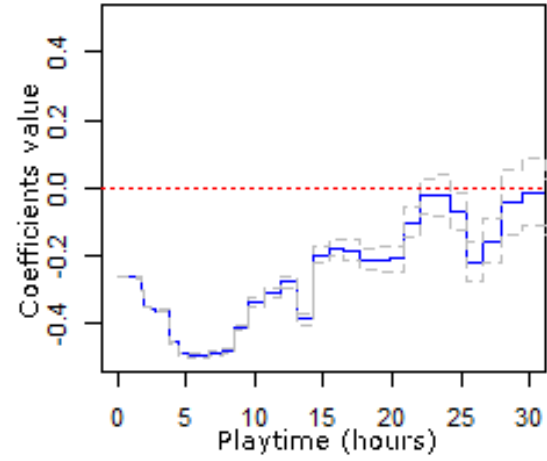
100 répétitions Monte-Carlo

# Application

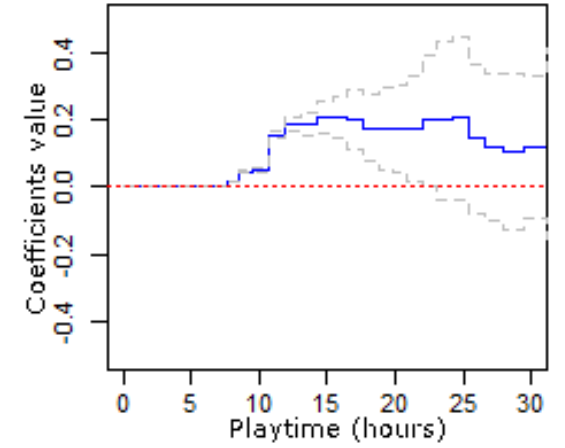
Les facteurs de rétention dans les jeux vidéo



### Rocket

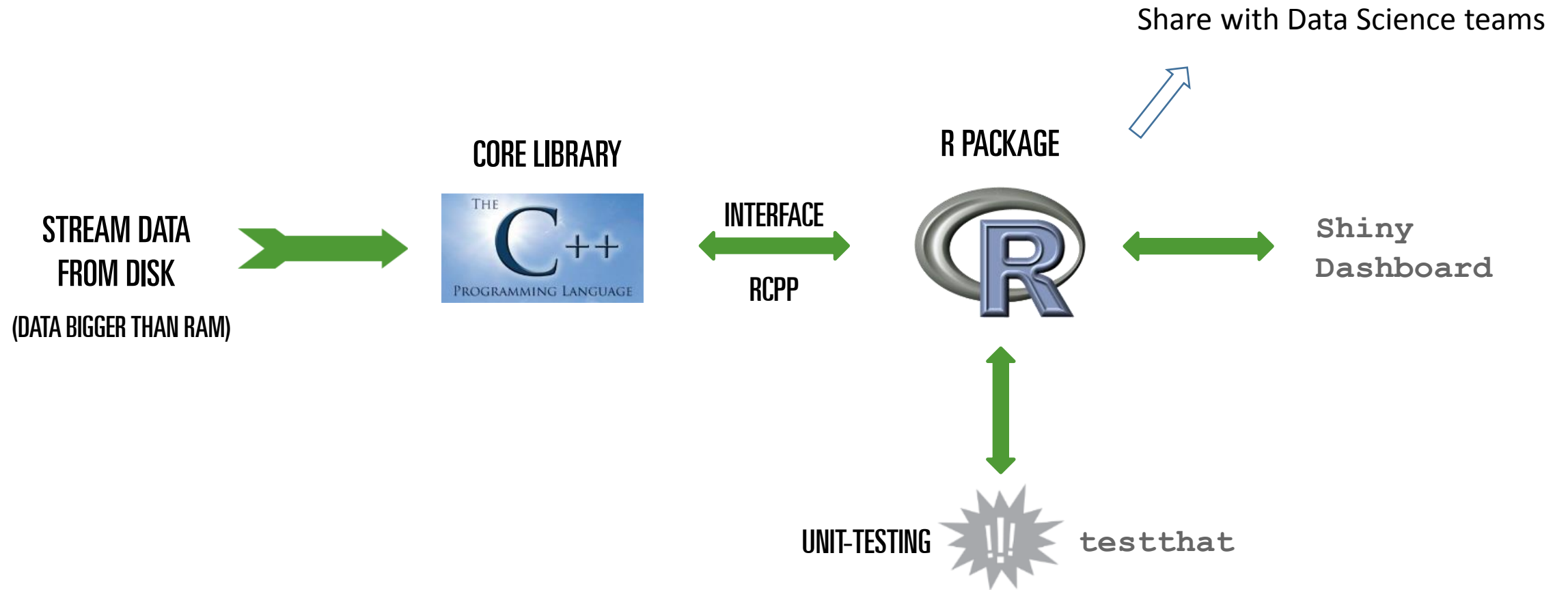


### Auto Crossbow



Merci !

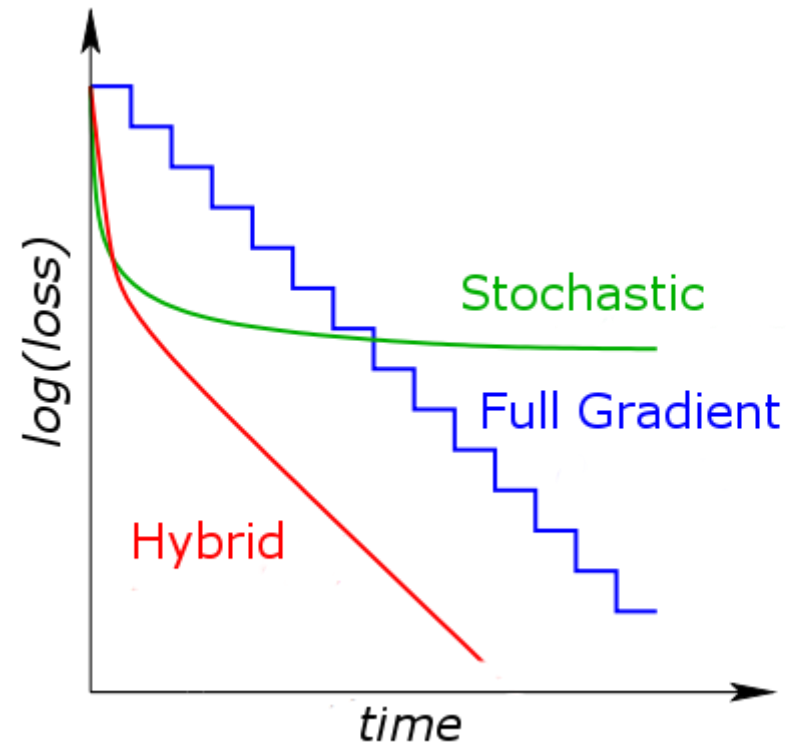
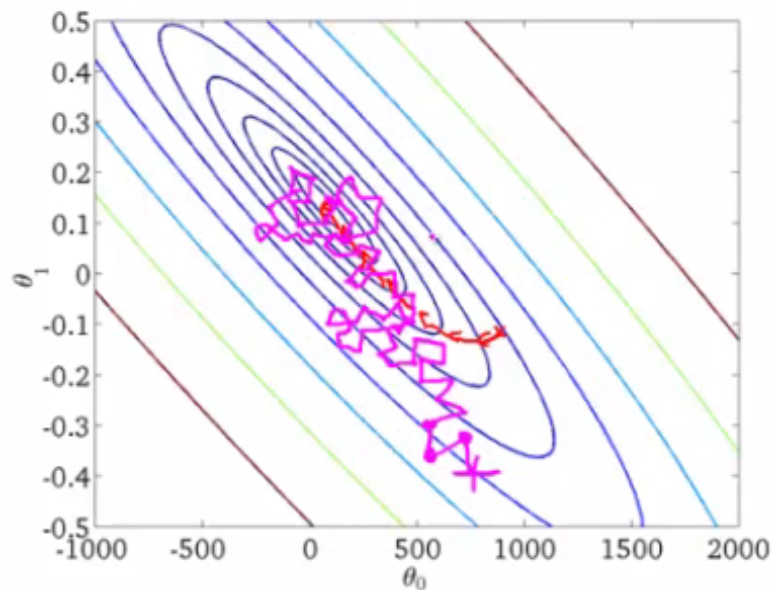
Bonus



# Stochastic Gradient Descent

$$\beta \leftarrow \beta - \mu \underbrace{\nabla_{\mathbf{i}, \beta}}_{\text{Gradient compute on one observation}}$$

Gradient compute on one observation





# Cox proportional likelihood

Cost function :

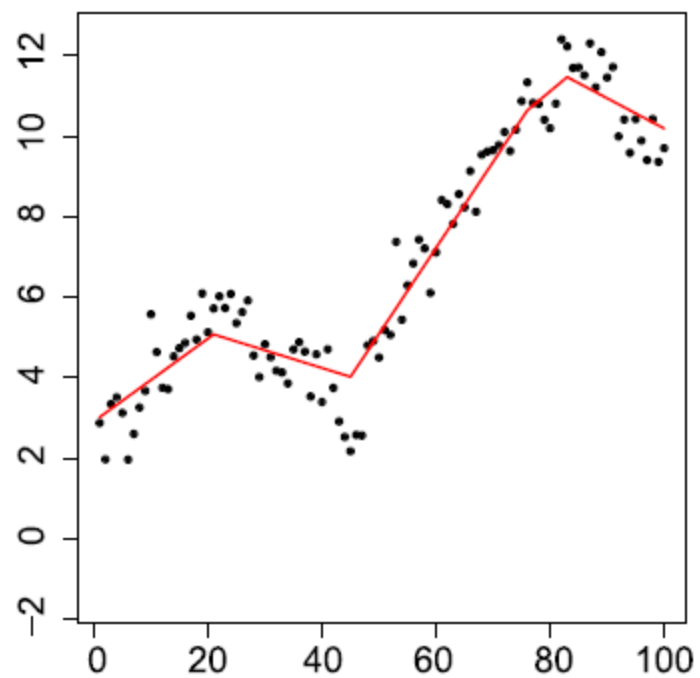
$$L(\beta) = \frac{1}{n} \prod_{i \in D} \frac{\exp(x_i^T \beta)}{\sum_{j \in R_i} \exp(x_j^T \beta)}$$

Individual gradients:

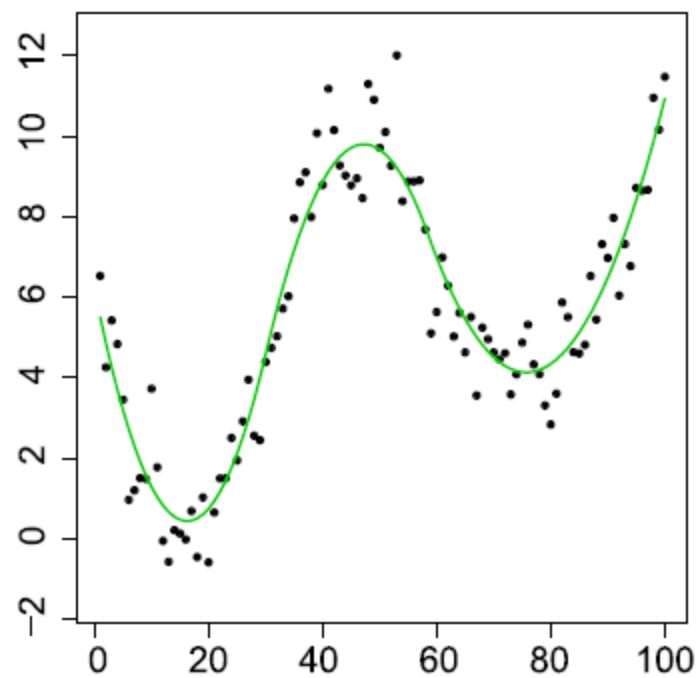
$$\Delta f_i(\beta) = -x_i + \sum_{j \in R_i} \frac{x_j \exp(x_j^T \beta)}{\sum_{k \in R_i} \exp(x_k^T \beta)} \quad \sim o(np)$$

It is  $o(p)$  for regression or logistic regression.

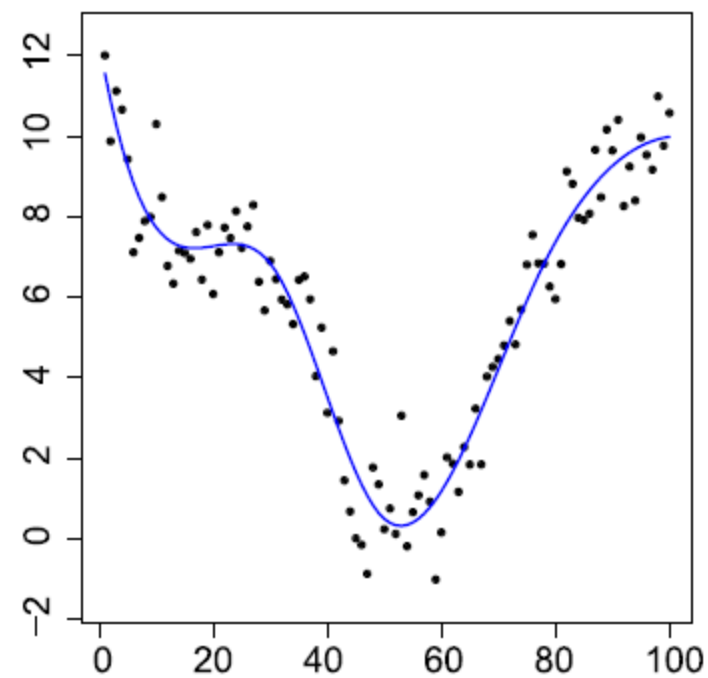
Computing individual gradient costs almost the same as computing full gradient !



(a) Linear

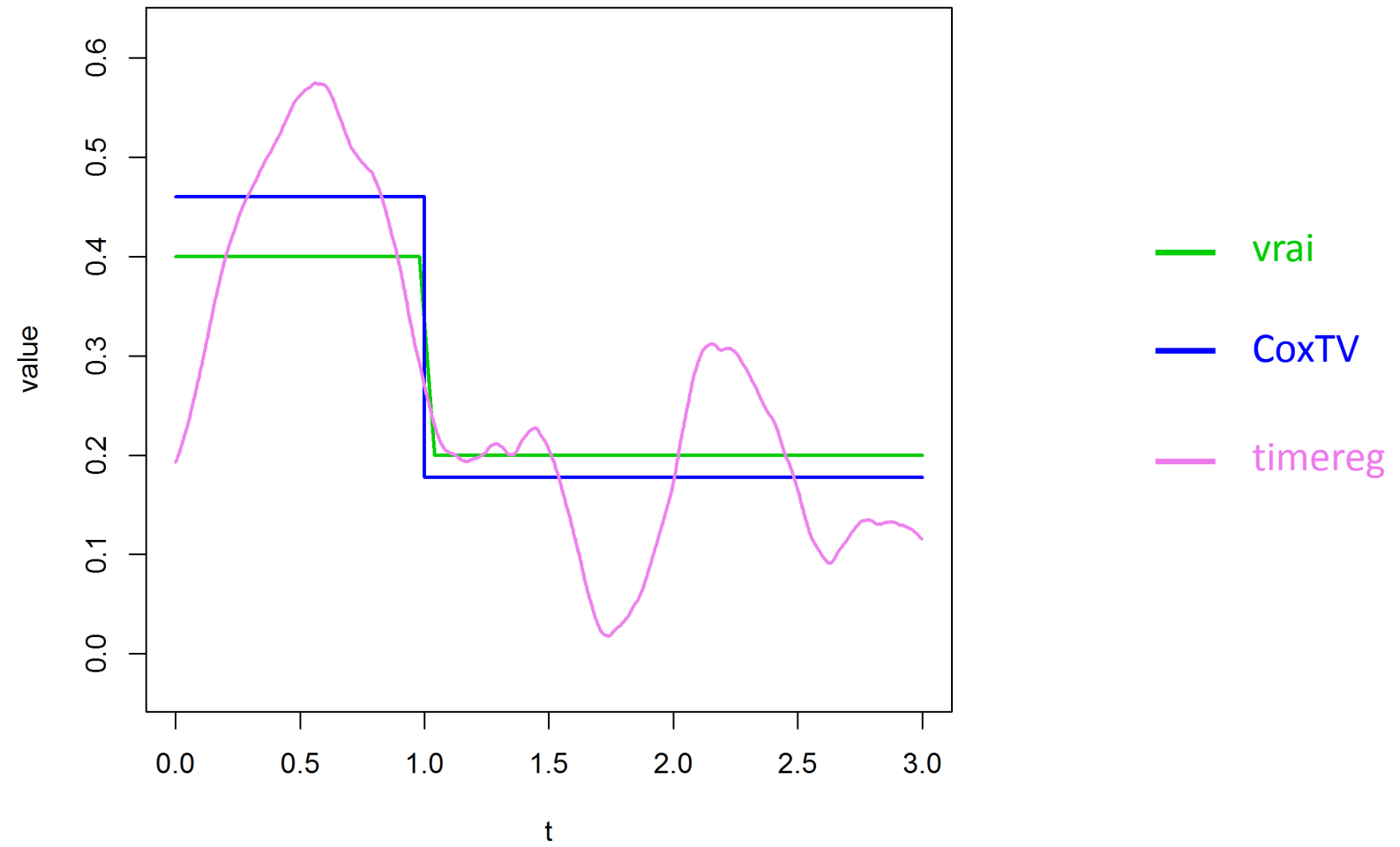


(b) Quadratic



(c) Cubic

**Beta 1**



$N = 1\,000$

$N_{\text{event}} \approx 3\,000$

$Nb_{\text{obs}} \approx 54\,000$

Nous introduisons la pénalité  $(\ell_1 + \ell_1)$ -variation totale avec poids défini par

$$\|\beta\|_{\text{gTV}, \hat{\gamma}} = \sum_{j=1}^p (\hat{\gamma}_{j,1} |\beta_{j,1}| + \sum_{l=2}^L \hat{\gamma}_{j,l} |\beta_{j,l} - \beta_{j,l-1}|)$$

pour tout  $\beta \in \mathbb{R}^{p \times L}$  avec  $\hat{\gamma} = (\hat{\gamma}_{1,\cdot}^\top, \dots, \hat{\gamma}_{p,\cdot}^\top)^\top$ , tel que  $\hat{\gamma}_{j,\cdot} \in \mathbb{R}_+^L$  pour tout  $j = 1, \dots, p$ , donné par

$$\hat{\gamma}_{j,l} \approx \sqrt{\frac{L \log(pL)}{n}} \hat{V}_{j,l}, \text{ avec } \hat{V}_{j,l} = \frac{1}{n} \sum_{i=1}^n \int_{\cup_{u=l}^L I_u} (X_i^j(t))^2 dN_i(t).$$

$$\hat{\beta}^{\text{M}} = \operatorname{argmin}_{\beta \in \mathbb{R}^{p \times L}} \left\{ \ell_n^{\text{M}}(\beta) + \|\beta\|_{\text{gTV}, \hat{\gamma}} \right\}.$$

divergence de Kullback empirique associée au modèle de Cox (2)

$$\begin{aligned} K_n(\lambda_\star^{\text{M}}, \lambda_\beta^{\text{M}}) &= \frac{1}{n} \sum_{i=1}^n \int_0^\tau (\log \lambda_\star^{\text{M}}(t, X_i(t)) - \log \lambda_\beta^{\text{M}}(t, X_i(t))) \lambda_\star^{\text{M}}(t, X_i(t)) Y_i(t) dt \\ &\quad - \frac{1}{n} \sum_{i=1}^n \int_0^\tau (\lambda_\star^{\text{M}}(t, X_i(t)) - \lambda_\beta^{\text{M}}(t, X_i(t))) Y_i(t) dt. \end{aligned}$$

**Théorème 2.** *Pour  $x > 0$  fixé, l'estimateur  $\hat{\lambda}^M$  défini dans (4), vérifie avec une probabilité supérieur à  $1 - C_M e^{-x}$  ( $C_M > 0$ ),*

$$K_n(\lambda_\star^M, \hat{\lambda}^M) \leq \inf_{\beta \in \mathbb{R}^{p \times L}} \left( K_n(\lambda_\star^M, \lambda_\beta^M) + 2\|\beta\|_{gTV, \hat{\gamma}} \right). \quad (6)$$

Ces théorèmes admettent l'interprétation suivante : les risques théoriques de nos estimateurs sont bornés par les meilleurs risques atteignables sur l'ensemble des modèles basés sur les histogrammes ( $\Lambda^A$  et  $\Lambda^M$ ), plus un terme satisfaisant

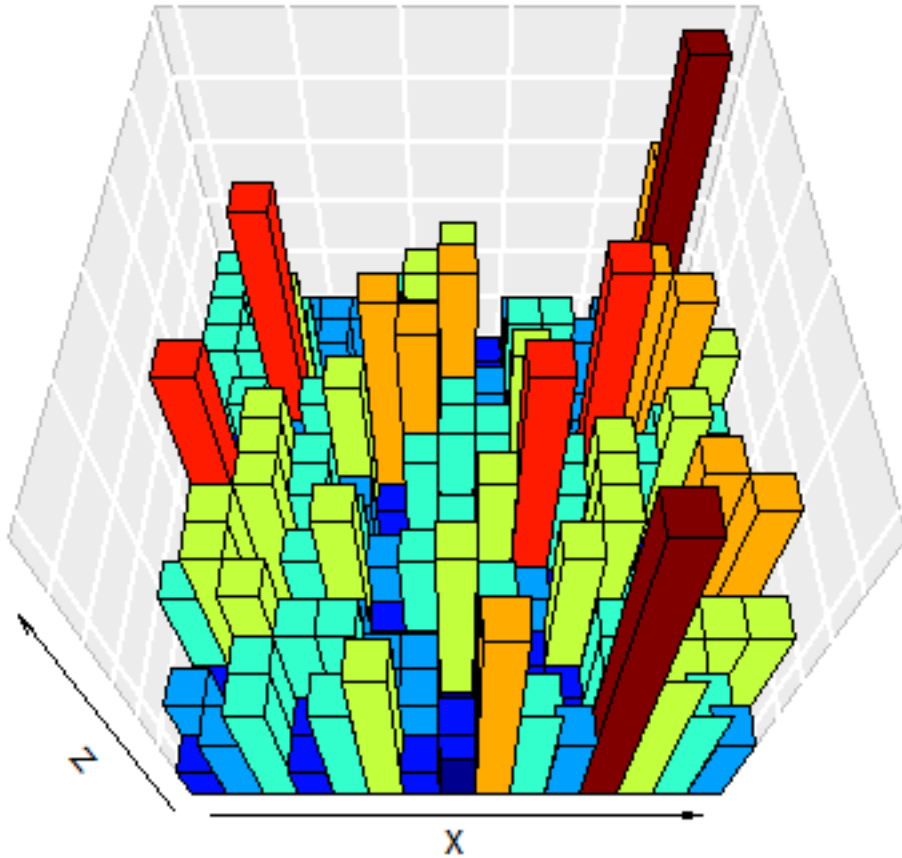
$$\|\beta\|_{gTV, \hat{\gamma}} \asymp \|\beta\|_{gTV} \max_{j=1, \dots, p} \max_{l=1, \dots, L} \sqrt{\frac{L \log pL}{n} \hat{V}_{j, \ell}}. \quad (7)$$

pour tout  $\beta \in \mathbb{R}^{p \times L}$ . La norme  $\|\cdot\|_{gTV}$  représente la  $(\ell_1 + \ell_1)$ -variation totale sans poids (*i.e.*  $\hat{\gamma}_{j, \ell} = 1$ ). Le terme dominant dans (7) est d'ordre  $\|\beta\|_{gTV} (L \log(pL)/n)^{1/2}$  (à un facteur de  $\log \log$  près).

$$n * p * t_{i,j}$$

$$n * p * t_i$$

$$t_i = \sum_{\text{temps observation}} \geq \max_j t_{i,j}$$



Id	event	start	stop	$x_1$	...	$x_p$
1001	0	0	1.2	1	...	-5
1001	0	1.2	2.0	3	...	-5
1001	0	2.0	5.3	5	...	-15
1001	1	5.3	5.8	8	...	22
1002	0	0	0.2	2	...	-1

1. Increase frequency for all covariates
2. Lose information

Or losing and transforming information if we discretize

## II

# Data

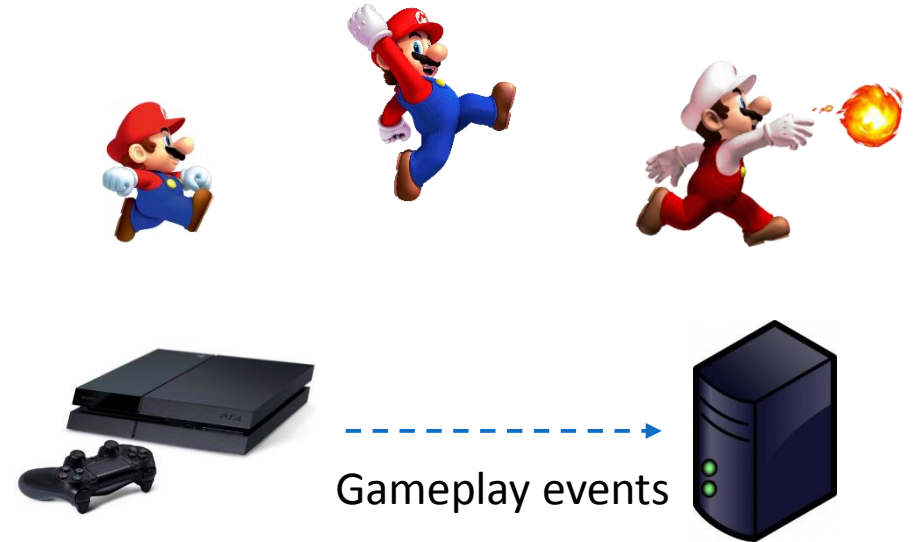
- We can track every change in the game



Levelling



Weapon usage



Datasets size :

millions of players \* thousands of observations per variable \* many variables

Player 1



Player 2

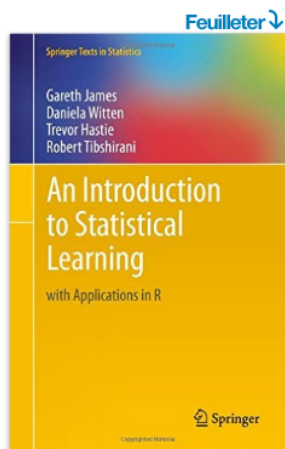


2h

30h Playtime







Feuilleter ↴

## An Introduction to Statistical Learning: With Applications in R (Anglais) R

de Gareth James ▾ (Auteur), Trevor Hastie ▾ (Auteur), Robert Tibshirani ▾ (Auteur), Daniela Witten ▾ (Auteur)

★★★★★ ▾ 128 Commentaires sur Amazon.com 🇺🇸 | Soyez la première personne à écrire un commentair

**#1 Meilleure vente** dans Mathematical & Statistical Software

▸ Voir les formats et éditions

Relié  
EUR 51,43

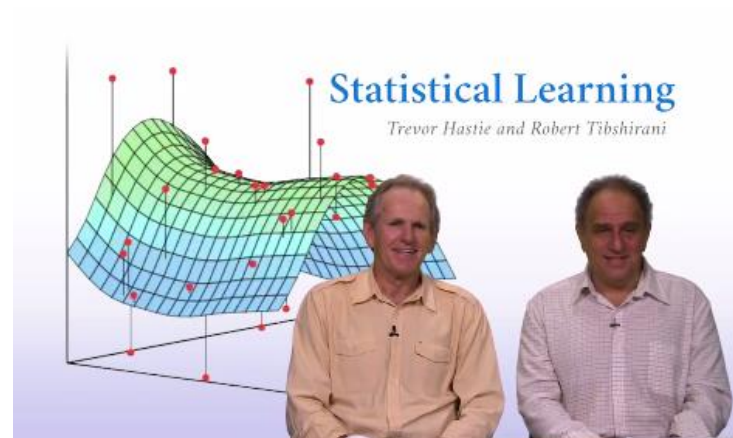
7 d'occasion à partir de EUR 51,63

17 neufs à partir de EUR 48,76

Voulez-vous le faire livrer aujourd'hui, mardi 10 jan.? Commandez-le dans les **3 h et 15 mins** et choisissez |

**Note:** Cet article est éligible à la livraison en points de collecte. [Détails](#)

**An Introduction to Statistical Learning** provides an accessible overview of the field of statistical learning, an essential toolset for making sense of the vast and complex data sets that have emerged in fields ranging from biology to finance to marketing to astrophysics in the past twenty years. This book



temps

★★★★★ wonderful but watch the movie 14 février 2014

Par I Teach Typing - **Publié sur Amazon.com**

Format: Relié | **Achat vérifié**

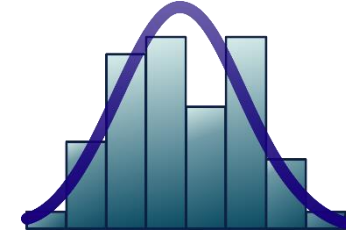
This is a wonderful book written by luminaries in the field. While it is not for casual consumption, it is a relatively approachable review of the state of the art for people who do not have the hardcore math needed for *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Second Edition (Springer Series in Statistics). This book is the text for the free Winter 2014 MOOC run out of Stanford called StatLearning (sorry Amazon will not allow me to include the website). Search for the class and you can watch Drs. Hastie and Tibshirani teach the material in this book.



↗ Playtime



↗ Lifetime



Survival Analysis



Repeated Events