

# Estimation régularisée du risque pour l'analyse age-period-cohort

---

V. Goepp<sup>†</sup>, G. Nuel<sup>‡</sup>, O. Bouaziz<sup>†</sup>

<sup>†</sup> : MAP5, Université Paris-Descartes

<sup>‡</sup> : LPMA, Université Pierre et Marie Curie

Séminaire de Statistiques, MAP5, 15 décembre 2017

# Plan de la présentation

- I Introduction
- II Modèles existants
- III Notre approche
- IV Résultats numériques : simulations
- V Résultats numériques : données réelles

# Introduction

## Présentation des données

Population: 91992 femmes adhérentes à la MGEN

Données récoltées par formulaire (2-3 ans)

Date calendaire  $\in [1990, 2010]$

L'évènement observé est l'apparition du cancer du sein

Objectif : estimer le risque d'avoir un cancer du sein à chaque instant

Difficultés :

- Pourcentage de cancers observés: 7%
- Date de naissance  $\in [1925, 1950]$  : population hétérogène

[1] F. Clavel-Chapelon et al, Cohort profile: the French E3N cohort study, *International journal of epidemiology*, 2014.

# Introduction

## Analyse de survie

- On veut estimer  $T$ , le temps passé avant l'apparition d'un évènement.
- On n'a pas accès à  $(T_i)_i$  mais à

$$Y_i = \min(T_i, C_i)$$

où  $C$  est une variable de censure avec  $C \perp\!\!\!\perp T$ .

- On connaît aussi  $\Delta_i = \mathbb{1}_{Y_i=T_i}$ .
- On estime le risque instantané:

$$\lambda(t) = \lim_{\delta t \rightarrow 0} \frac{\mathbb{P}(t \leq T \leq t + \delta t | T > t)}{\delta t}$$

# Introduction

## Diagramme de Lexis

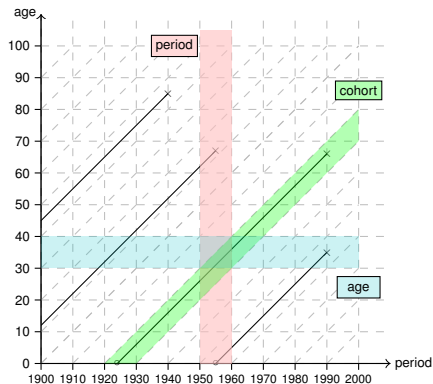


Diagramme de Lexis: Age-Period

# Introduction

## Diagramme de Lexis

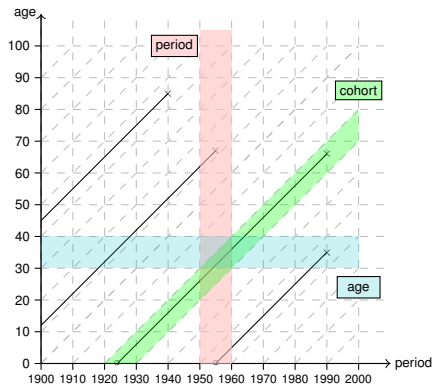


Diagramme de Lexis: Age-Period

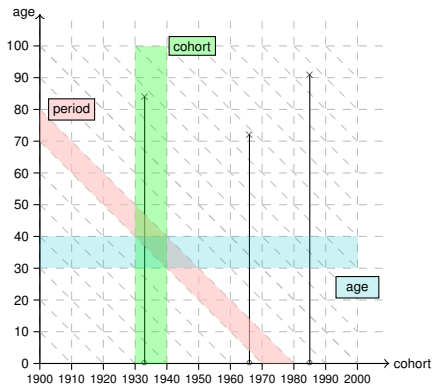


Diagramme Age-Cohort

Le risque instantané  $\lambda$  est discrétisé en  $J$  intervalles d'âge et  $K$  intervalles de cohorte :

$$\lambda(\text{age}, \text{cohorte}) = \sum_{j=1}^J \sum_{k=1}^K \lambda_{j,k} \mathbb{1}_{[c_{j-1}, c_j) \times [d_{k-1}, d_k)}(\text{age}, \text{cohorte})$$

Objectif : estimer  $\lambda_{j,k}$

# Analyse age-period-cohort

On veut modéliser l'effet de l'âge, la cohorte et la période.

- effet de l'âge : ménopause
- effet de la cohorte : biberon cancérigène
- effet de la période : accident nucléaire

On définit un vecteur de paramètres par effet:  $\alpha$ ,  $\beta$  et  $\gamma$



# Modèles existants

(i) Dans le modèle AGE-COHORT, on suppose

$$\log \lambda_{j,k} = \alpha_j + \beta_k.$$

- $J + K - 1$  paramètres pour  $JK$  variables: régularisation
- Fort *a priori* sur  $\lambda$

(ii) Dans le modèle AGE-PERIOD-COHORT, on suppose

$$\log \lambda_{j,k} = \alpha_j + \beta_k + \gamma_{j+k-1}.$$

- Non identifiable : on peut soit
  - estimer  $\Delta^2 \alpha$ ,  $\Delta^2 \beta$  et  $\Delta^2 \gamma$ .
  - rajouter une contrainte.

[2] B. Carstensen, Age–period–cohort models for the Lexis diagram, *Statistics in medicine*, 2007.

# Estimateur du maximum de vraisemblance

On appelle :

- $O_{j,k}$  : nombre d'évènements dans le  $(j, k)$ -ième rectangle
- $R_{j,k}$  : temps à risque dans le  $(j, k)$ -ième rectangle

La log vraisemblance négative s'écrit

$$\ell_n(\lambda) = \sum_{j=1}^J \sum_{k=1}^K \lambda_{j,k} R_{j,k} - O_{j,k} \log(\lambda_{j,k}) .$$

L'estimateur du maximum de vraisemblance est :

$$\lambda_{j,k}^{\text{mle}} = \frac{O_{j,k}}{R_{j,k}}$$

# Estimateur du maximum de vraisemblance

On appelle :

- $O_{j,k}$  : nombre d'évènements dans le  $(j, k)$ -ième rectangle
- $R_{j,k}$  : temps à risque dans le  $(j, k)$ -ième rectangle

La log vraisemblance négative s'écrit

$$\ell_n(\lambda) = \sum_{j=1}^J \sum_{k=1}^K \lambda_{j,k} R_{j,k} - O_{j,k} \log(\lambda_{j,k}) .$$

L'estimateur du maximum de vraisemblance est :

$$\lambda_{j,k}^{\text{mle}} = \frac{O_{j,k}}{R_{j,k}} \rightarrow \text{overfitting}$$

# Notre approche : vraisemblance pénalisée

Aucun *a priori* :

$$\log \lambda_{j,k} = \eta_{j,k},$$

# Notre approche : vraisemblance pénalisée

Aucun *a priori* :

$$\log \lambda_{j,k} = \eta_{j,k},$$

Mais l'estimation de  $\eta$  est faite par vraisemblance **pénalisée**:

$$\ell_n^{\text{pen}}(\eta) = \underbrace{\ell_n(\eta)}_{\text{attache aux données}}$$

# Notre approche : vraisemblance pénalisée

Aucun *a priori* :

$$\log \lambda_{j,k} = \eta_{j,k},$$

Mais l'estimation de  $\eta$  est faite par vraisemblance **pénalisée**:

$$\ell_n^{\text{pen}}(\eta) = \underbrace{\ell_n(\eta)}_{\text{attache aux données}} + \underbrace{\frac{\text{pen}}{2} \sum_{j,k} v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2 + w_{j,k} (\eta_{j,k+1} - \eta_{j,k})^2}_{\text{régularisation}},$$

# Notre approche : vraisemblance pénalisée

Aucun *a priori* :

$$\log \lambda_{j,k} = \eta_{j,k},$$

Mais l'estimation de  $\eta$  est faite par vraisemblance **pénalisée**:

$$\ell_n^{\text{pen}}(\eta) = \underbrace{\ell_n(\eta)}_{\text{attache aux données}} + \underbrace{\frac{\text{pen}}{2} \sum_{j,k} v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2 + w_{j,k} (\eta_{j,k+1} - \eta_{j,k})^2}_{\text{régularisation}},$$

$v$  et  $w$  sont des poids,

pen est une constante de régularisation.

# Deux types de régularisation

(i) Régularisation  $L_2$  (Ridge) avec  $\mathbf{v} = \mathbf{w} = \mathbf{1}$



# Deux types de régularisation

- (i) Régularisation  $L_2$  (Ridge) avec  $\mathbf{v} = \mathbf{w} = \mathbf{1}$
- (ii) Régularisation  $L_0$  avec la procédure itérative **adaptive ridge**.  
Les poids sont adaptés itérativement :

$$\begin{cases} v_{j,k} = \left( \left( \eta_{j+1,k} - \eta_{j,k} \right)^2 + \varepsilon^2 \right)^{-1} \\ w_{j,k} = \left( \left( \eta_{j,k} - \eta_{j,k-1} \right)^2 + \varepsilon^2 \right)^{-1}, \end{cases}$$

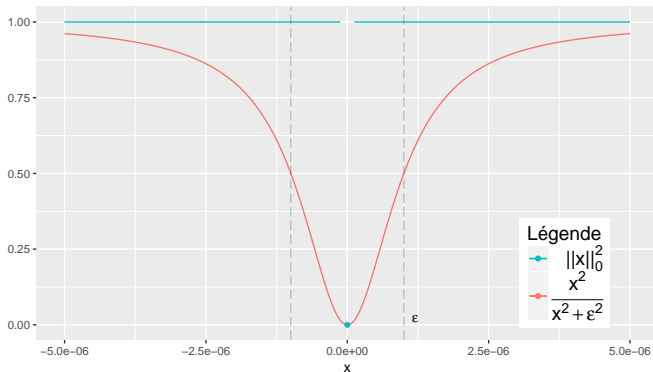
avec  $\varepsilon \ll 1$ .

- [3] F. Frommlet and G. Nuel, An Adaptive Ridge Procedure for  $L_0$  Regularization, *Public Library of Science*, 2016.

# Approximation de la norme $L_0$

Lorsque  $\varepsilon \ll 1$ :

$$v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2 \simeq \|\eta_{j+1,k} - \eta_{j,k}\|_0^2 = \begin{cases} 0 & \text{si } \eta_{j+1,k} = \eta_{j,k} \\ 1 & \text{si } \eta_{j+1,k} \neq \eta_{j,k} \end{cases}$$



# La procédure *Adaptive Ridge*

**procedure** ADAPTIVE-RIDGE( $\mathbf{O}$ ,  $\mathbf{R}$ , pen)

**end procedure**

# La procédure *Adaptive Ridge*

**procedure** ADAPTIVE-RIDGE(***O***, ***R***, pen)

**$\eta$**   $\leftarrow$  **0**

**$v$**   $\leftarrow$  **1**

**$w$**   $\leftarrow$  **1**

**end procedure**

# La procédure *Adaptive Ridge*

**procedure** ADAPTIVE-RIDGE( $\mathbf{O}, \mathbf{R}, \text{pen}$ )

$\boldsymbol{\eta} \leftarrow \mathbf{0}$

$\mathbf{v} \leftarrow \mathbf{1}$

$\mathbf{w} \leftarrow \mathbf{1}$

**while** not converge **do**

$\boldsymbol{\eta}^{\text{new}} \leftarrow \text{NEWTON-RAPHSON}(\mathbf{O}, \mathbf{R}, \text{pen}, \mathbf{v}, \mathbf{w})$

$\mathbf{v}_{j,k}^{\text{new}} \leftarrow \left( (\eta_{j+1,k}^{\text{new}} - \eta_{j,k}^{\text{new}})^2 + \varepsilon^2 \right)^{-1}$

$\mathbf{w}_{j,k}^{\text{new}} \leftarrow \left( (\eta_{j,k}^{\text{new}} - \eta_{j,k-1}^{\text{new}})^2 + \varepsilon^2 \right)^{-1}$

$\boldsymbol{\eta} \leftarrow \boldsymbol{\eta}^{\text{new}}$

$\mathbf{v} \leftarrow \mathbf{v}^{\text{new}}$

$\mathbf{w} \leftarrow \mathbf{w}^{\text{new}}$

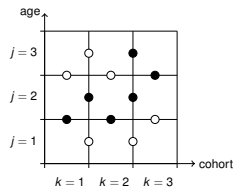
**end while**

**end procedure**

# La procédure *Adaptive Ridge*

```
procedure ADAPTIVE-RIDGE( $\mathbf{O}, \mathbf{R}, \text{pen}$ )  
   $\eta \leftarrow \mathbf{0}$   
   $\mathbf{v} \leftarrow \mathbf{1}$   
   $\mathbf{w} \leftarrow \mathbf{1}$   
  while not converge do  
     $\eta^{\text{new}} \leftarrow \text{NEWTON-RAPHSON}(\mathbf{O}, \mathbf{R}, \text{pen}, \mathbf{v}, \mathbf{w})$   
     $\mathbf{v}_{j,k}^{\text{new}} \leftarrow \left( (\eta_{j+1,k}^{\text{new}} - \eta_{j,k}^{\text{new}})^2 + \varepsilon^2 \right)^{-1}$   
     $\mathbf{w}_{j,k}^{\text{new}} \leftarrow \left( (\eta_{j,k}^{\text{new}} - \eta_{j,k-1}^{\text{new}})^2 + \varepsilon^2 \right)^{-1}$   
     $\eta \leftarrow \eta^{\text{new}}$   
     $\mathbf{v} \leftarrow \mathbf{v}^{\text{new}}$   
     $\mathbf{w} \leftarrow \mathbf{w}^{\text{new}}$   
  end while  
  Compute ( $\mathbf{O}^{\text{sel}}, \mathbf{R}^{\text{sel}}$ ) from ( $\eta^{\text{new}}, \mathbf{v}^{\text{new}}, \mathbf{w}^{\text{new}}$ )  
   $\eta^{\text{mle}} \leftarrow \mathbf{O}^{\text{sel}} / \mathbf{R}^{\text{sel}}$   
  return  $\eta^{\text{mle}}$   
end procedure
```

# Adaptive Ridge permet de sélectionner un modèle

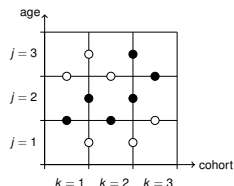


(a) Représentation de

$$v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2$$

et  $w_{j,k} (\eta_{j,k+1} - \eta_{j,k})^2$

# Adaptive Ridge permet de sélectionner un modèle

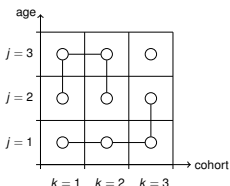


(a) Représentation de

$$v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2$$

et

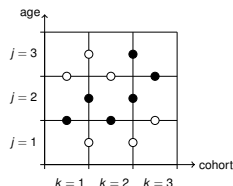
$$w_{j,k} (\eta_{j,k+1} - \eta_{j,k})^2$$



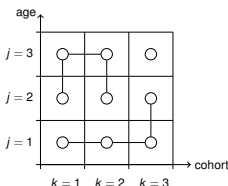
(b) Graphe correspondant



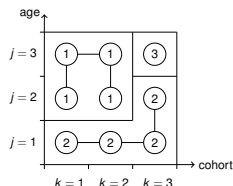
# Adaptive Ridge permet de sélectionner un modèle



(a) Représentation de  
 $v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2$   
 et  $w_{j,k} (\eta_{j,k+1} - \eta_{j,k})^2$



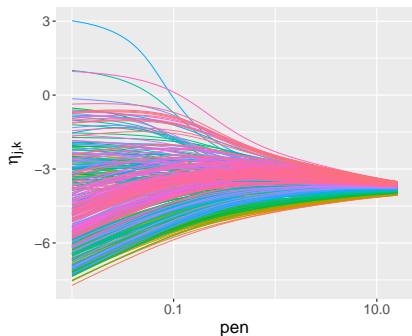
(b) Graphe correspondant



(c) Segmentation selon les  
 composantes connexes

# Comparaison des deux régularisations

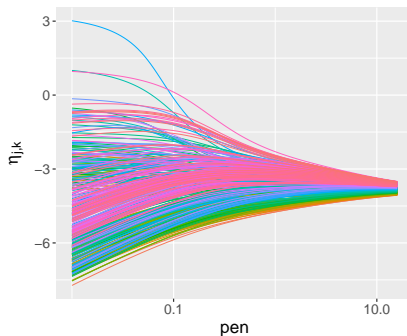
$\text{pen} \rightarrow 0$  :  $\hat{\lambda} \rightarrow \hat{\lambda}^{\text{mle}}$   
 $\text{pen} \rightarrow \infty$  :  $\hat{\lambda}$  constant



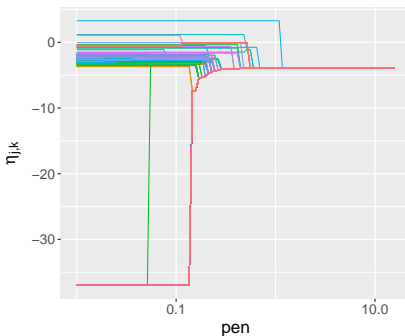
Régularisation  $L_2$  : à chaque pen  
correspond un estimateur

# Comparaison des deux régularisations

$\text{pen} \rightarrow 0$  :  $\hat{\lambda} \rightarrow \hat{\lambda}^{\text{mle}}$   
 $\text{pen} \rightarrow \infty$  :  $\hat{\lambda}$  constant



Régularisation  $L_2$  : à chaque  $\text{pen}$  correspond un estimateur



Régularisation  $L_0$  : à chaque  $\text{pen}$  correspond un *modèle*

# Choix du modèle pour *Adaptive Ridge*

## Critères bayésiens

- Problème: choisir entre  $M$  modèles  $\mathcal{M}_1, \dots, \mathcal{M}_M$  de dimensions  $q_1, \dots, q_M$ .
- Solution: maximiser  $\mathbb{P}(\mathcal{M}_m | \mathbf{R}, \mathbf{O}) \propto \mathbb{P}(\mathbf{R}, \mathbf{O} | \mathcal{M}_m) \pi(\mathcal{M}_m)$ .
- Par approximation :  
$$-2 \log (\mathbb{P}(\mathcal{M}_m | \mathbf{R}, \mathbf{O})) = 2\ell_n(\hat{\eta}_m) + q_m \log n - 2 \log \pi(\mathcal{M}_m) + \mathcal{O}_{\mathbb{P}}(1)$$
- Comment choisir la distribution *a priori*  $\pi(\mathcal{M}_m)$  ?

# Choix du modèle pour *Adaptive Ridge*

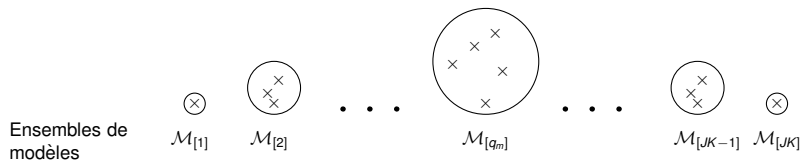
BIC:  $\pi(\mathcal{M}_m) = 1$

Tous les  $\mathcal{M}_m$  sont équiprobables

# Choix du modèle pour *Adaptive Ridge*

BIC:  $\pi(\mathcal{M}_m) = 1$

Tous les  $\mathcal{M}_m$  sont équiprobables

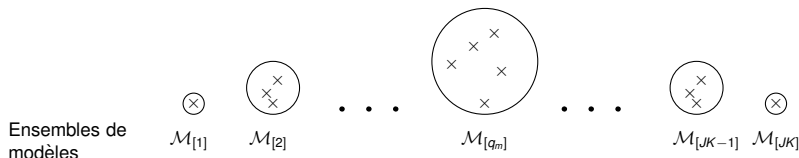


$\mathcal{M}_{[q_m]}$  est l'ensemble des modèles de dimension  $q_m$

# Choix du modèle pour *Adaptive Ridge*

BIC:  $\pi(\mathcal{M}_m) = 1$   
Tous les  $\mathcal{M}_m$  sont équiprobables

EBIC<sub>0</sub>:  $\mathbb{P}(\mathcal{M}_m \in \mathcal{M}_{[q_m]}) = 1$   
Tous les  $\mathcal{M}_{[q_m]}$  sont équiprobables



$\mathcal{M}_{[q_m]}$  est l'ensemble des modèles de dimension  $q_m$

[4] J. Chen and Z. Chen, Extended Bayesian information criteria for model selection with large model spaces, *Biometrika*, 2008.

# Choix du modèle pour *Adaptive Ridge*

## Critères utilisés

On compare différents critères de sélection :

- (i)  $\text{BIC}(m) = 2\ell_n(\hat{\eta}_m) + q_m \log n$
- (ii)  $\text{EBIC}_0(m) = 2\ell_n(\hat{\eta}_m) + q_m \log n - 2 \log \binom{JK}{q_m}$
- (iii)  $\text{AIC}(m) = 2\ell_n(\hat{\eta}_m) + 2q_m$
- (iv) K-fold Cross validation (CV)



# Illustration sur données simulées

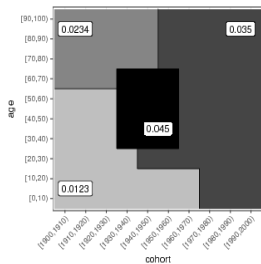
On simule des données selon :

- $\lambda$  constant par morceaux
- $\lambda$  lisse

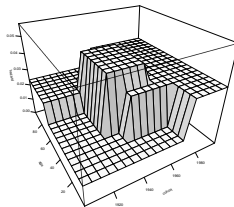
On compare :

- Modèle AGE-COHORT :  $\log \lambda_{j,k} = \alpha_j + \beta_k$
- Régularisation  $L_2$  avec CV
- Régularisation  $L_0$  avec AIC, BIC, EBIC<sub>0</sub> et CV.

# Simulations : cas n°1

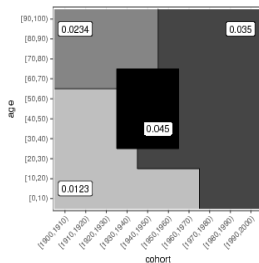


Vrai  $\lambda$

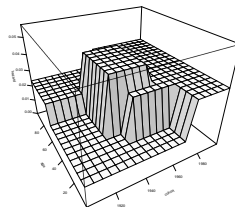


Vrai  $\lambda$

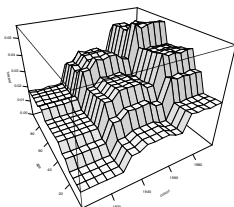
# Simulations : cas n°1



Vrai  $\lambda$

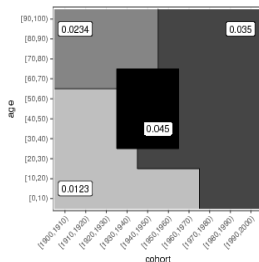


Vrai  $\lambda$

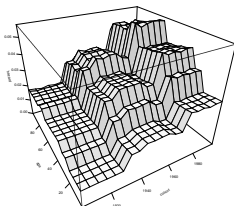


Modèle AGE-COHORT

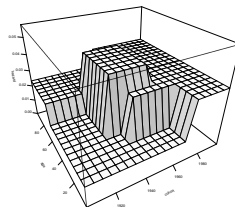
# Simulations : cas n°1



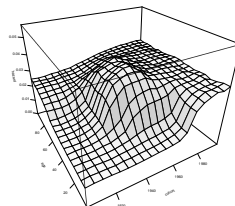
Vrai  $\lambda$



Modèle AGE-COHORT

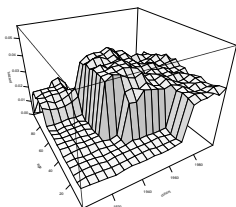


Vrai  $\lambda$



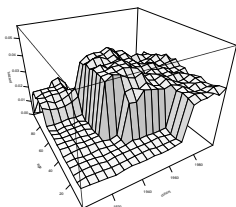
Régularisation  $L_2$  avec CV

# Simulations : cas n°1

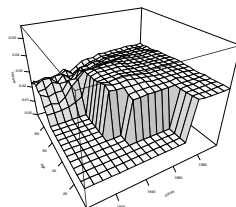


Régularisation  $L_0$  avec AIC

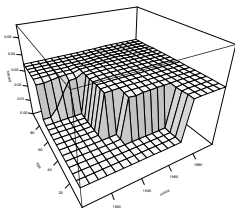
# Simulations : cas n°1



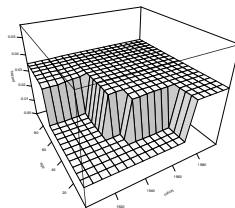
Régularisation  $L_0$  avec AIC



Régularisation  $L_0$  avec BIC

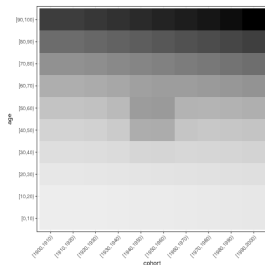


Régularisation  $L_0$  avec  $EBIC_0$

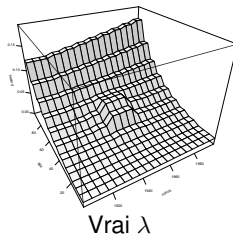


Régularisation  $L_0$  avec CV

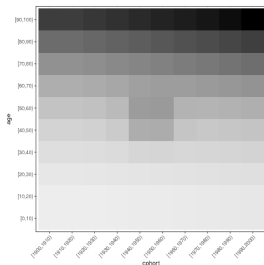
# Simulations : cas n°2



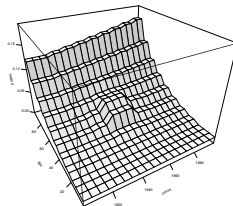
Vrai  $\lambda$



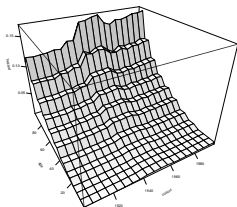
# Simulations : cas n°2



Vrai  $\lambda$



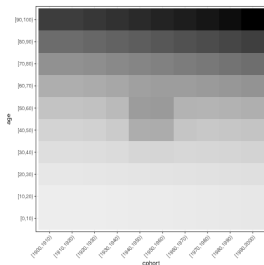
Vrai  $\lambda$



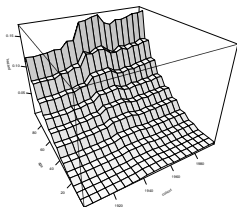
Modèle AGE-COHORT



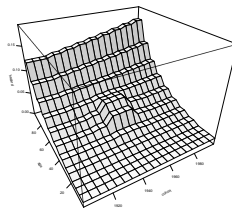
# Simulations : cas n°2



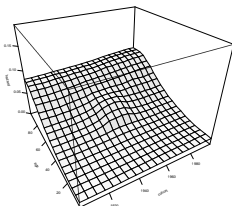
Vrai  $\lambda$



Modèle AGE-COHORT

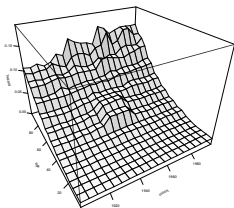


Vrai  $\lambda$



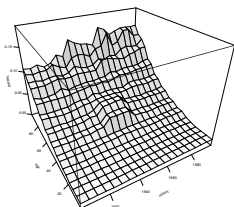
Régularisation  $L_2$  avec CV

# Simulations : cas n°2

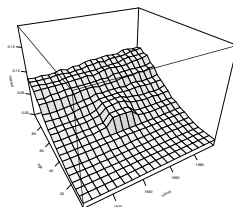


Régularisation  $L_0$  avec AIC

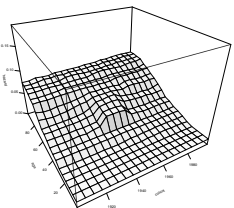
# Simulations : cas n°2



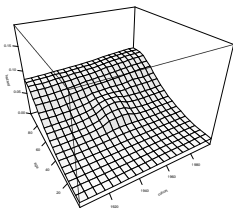
Régularisation  $L_0$  avec AIC



Régularisation  $L_0$  avec BIC



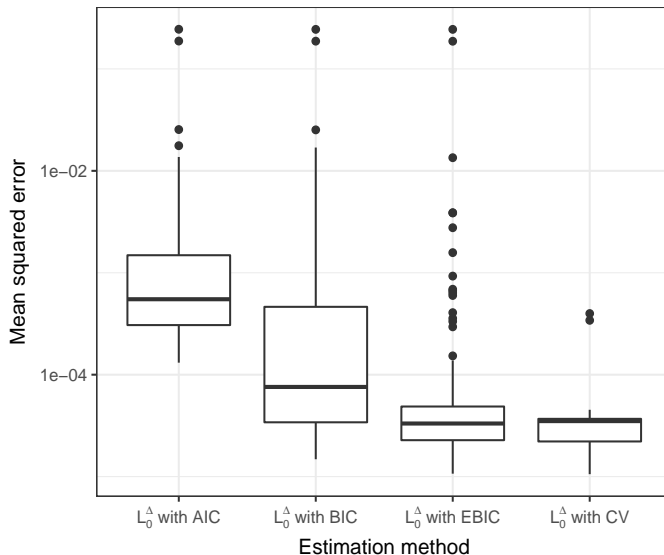
Régularisation  $L_0$  avec  $EBIC_0$



Régularisation  $L_0$  avec CV

# Simulations : comparaison quantitative

Erreur quadratique moyenne, pour  $\lambda$  constant par morceaux.

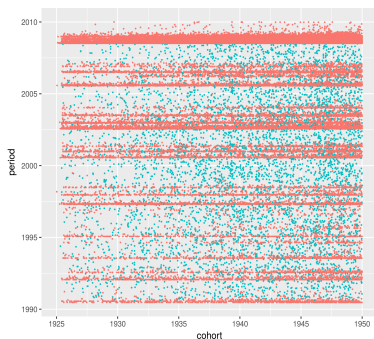


# Application : données réelles

## Présentation des données

Cohorte  $\in [1925, 1950]$

Période  $\in [1990, 2010]$



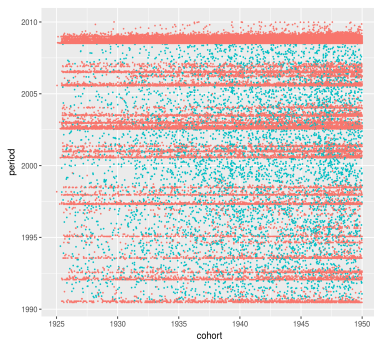
Plan période-cohorte

# Application : données réelles

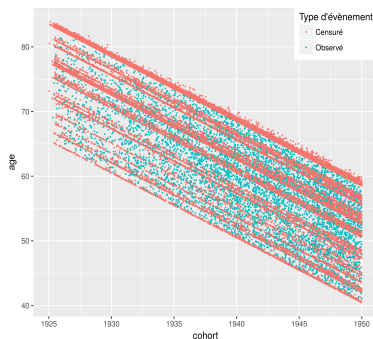
## Présentation des données

Cohorte  $\in [1925, 1950]$

Période  $\in [1990, 2010]$



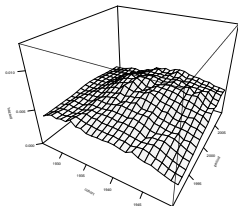
Plan période-cohorte



Plan âge-cohorte

# Application : données de l'étude E3N

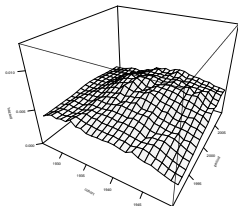
## Résultats



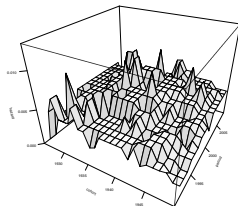
Régularisation  $L_2$  avec CV

# Application : données de l'étude E3N

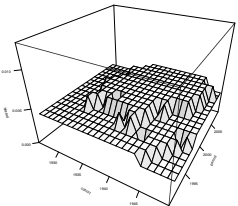
## Résultats



Régularisation  $L_2$  avec CV



Régularisation  $L_0$  avec AIC

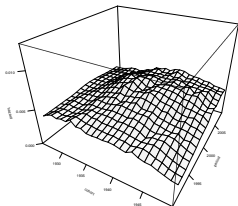


Régularisation  $L_0$  avec EBIC<sub>0</sub>

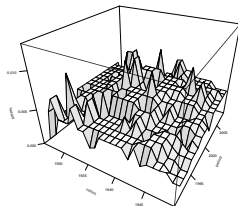


# Application : données de l'étude E3N

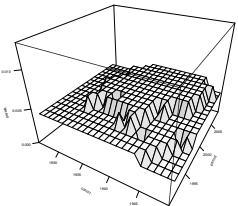
## Résultats



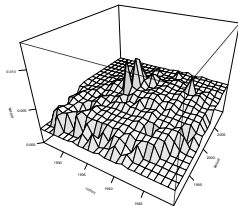
Régularisation  $L_2$  avec CV



Régularisation  $L_0$  avec AIC



Régularisation  $L_0$  avec  $EBIC_0$



Régularisation  $L_0$  avec  $EBIC_0$ :  
bootstrap

# Conclusion et perspectives

- Segmentation du risque instantané
- EBIC<sub>0</sub> plus performant que les autres critères
- Amélioration possible : différences d'ordre supérieur
- Le modèle peut s'étendre :

$$\log \lambda_{j,k} = \mu + \alpha_j + \beta_k + \delta_{j,k},$$

avec régularisation de  $\delta_{j,k}$

Merci