

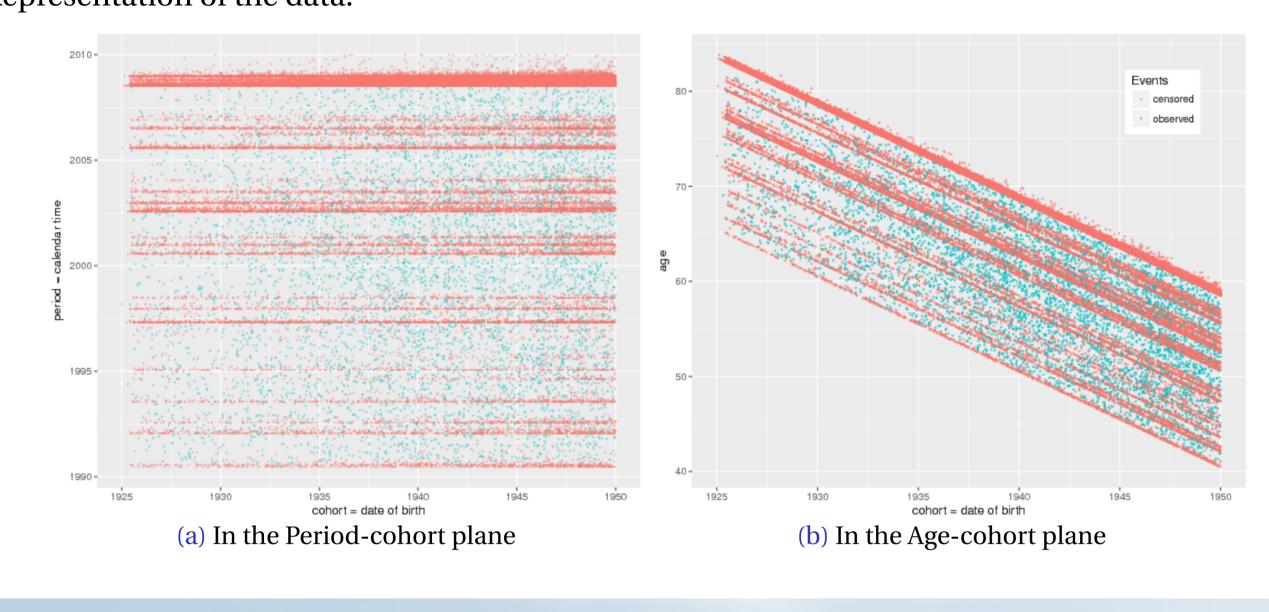


Presentation of the problem

The **E3N Cohort Study** ^[1]:

- Epigemiological study focused on the link between cancer and nutrition (part of EPIC).
- Population: ~ 100000 women (mostly teachers).
- Medical data gathered every 2-3 years using questionnaires (9 questionnaires in total). Diet, physical exercice and medical treatments are monitored.
- Blood and saliva samples are also gathered for som participants (not used here).
- The event of interest is the occurrence of breast cancer
- These occurrences are spread over the period [1990,2010]

Objective: estimate the hazard rate of breast cancer occurrence Representation of the data:



Right-Censoring

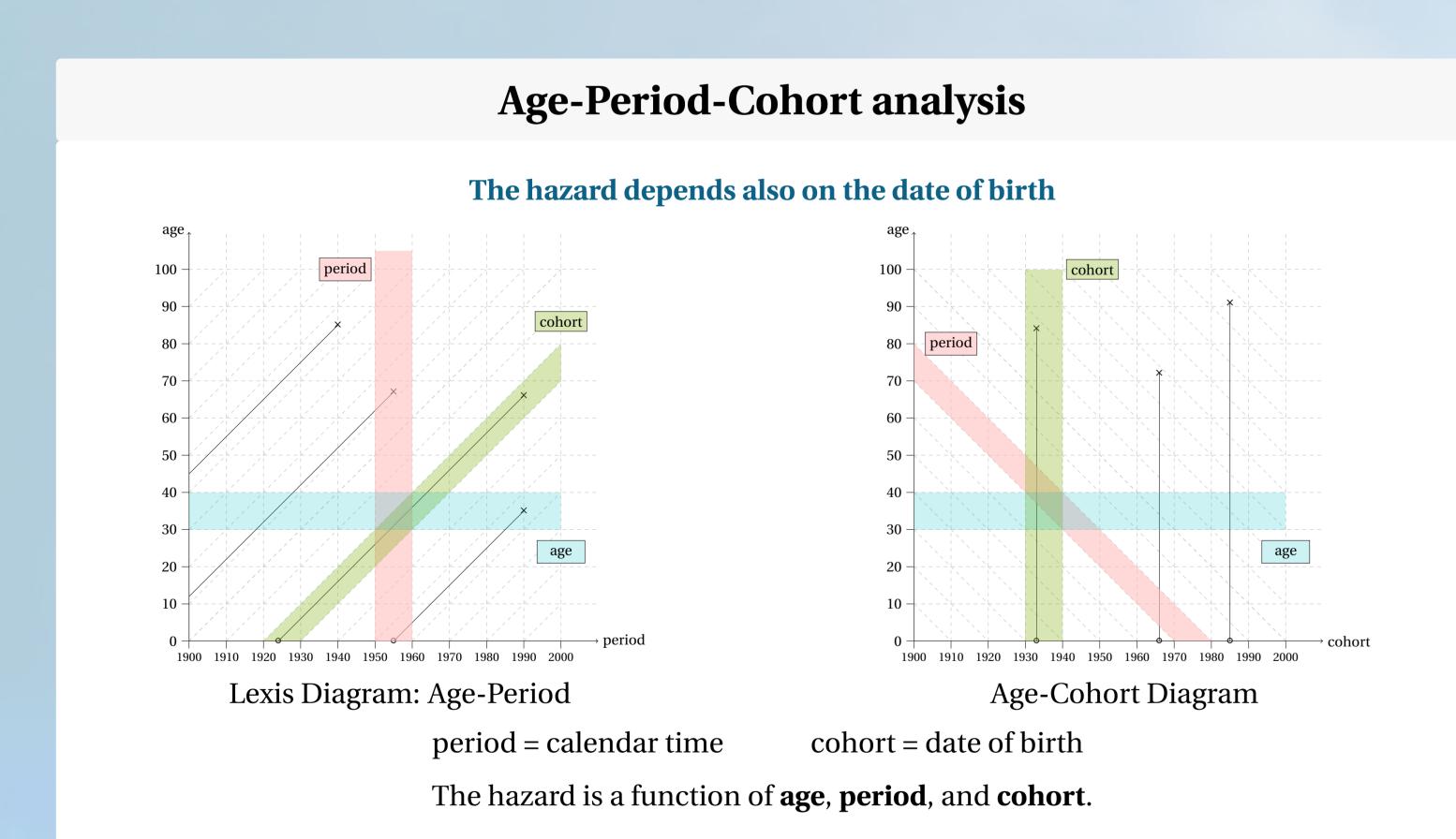
The age of cancer occurrence is observed for only 7% of the individuals.

- T_i is the age of cancer onset.
- We do not observe $(T_i)_i$ but

 $Y_i = \min(T_i, C_i)$ where *C* is a censoring r.v. independent from *T*.

- We observe $\Delta_i = 1_{T_i = Y_i}$.
- We infer the instantenous hazard rate

$$\lambda(t) = \lim_{dt \to 0} \frac{\mathbb{P}\left(t \le T \le t + dt | T \ge t\right)}{dt}$$



Existing Models in age-period-cohort Analysis

- In the literature: we infer lpha , eta , and γ , parameters of the age, cohort and period effetcs.
- In the AGE-COHORT Model, it is assumed that

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

J + *K* – 1 parameters for *JK* variables: regularizing
Strong *a prior* on λ.

► In the AGE-PERIOD-COHORT Model, it is assumed that

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

- Regularizing
 Strong *a priori* on λ
- Non identifiable



Regularized hazard estimation for age-period-cohort analysis

Vivien Goepp^{1,*}, Grégory Nuel², and Olivier Bouaziz²

Statistical Methods for Post Genomic Data, January 11-12 2018, Montpellier, France

New Approach: Penalized Likelihood

Reparametrization:

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

The unpenalized negative log-likelihood
$$\ell_n$$
 takes the form

$$\mathcal{C}_n(\boldsymbol{\eta}) = \sum_{j=1}^J \sum_{k=1}^K \exp(\eta_{j,k}) R_{j,k} - \eta_{j,k} O_{j,k}$$

where

- $O_{j,k}$ = number of observed events in the (j, k)-th rectangle
- $R_{j,k}$ = total time at risk in the (j, k)-th rectangle The MLE is explicit:

 $\hat{\eta}_{j,k}^{\text{mle}} = \log\left(\frac{O_{j,k}}{R_{j,k}}\right) \rightarrow \text{overfitting.}$

Our model has no *a priori*: But the inference is made by minimizing the **penalized likelihood** ^[2]

$$\ell_n^{\text{pen}}(\boldsymbol{\eta}) = \ell_n(\boldsymbol{\eta}) + \underbrace{\frac{\text{pen}}{2} \sum_{j,k} v_{j,k} (\eta_{j+1,k} - \eta_{j,k})^2 + w_{j,k} (\eta_{j,k+1} - \eta_{j,k})^2}_{\text{goodness}}$$
regularization

► *v* et *w* are weights

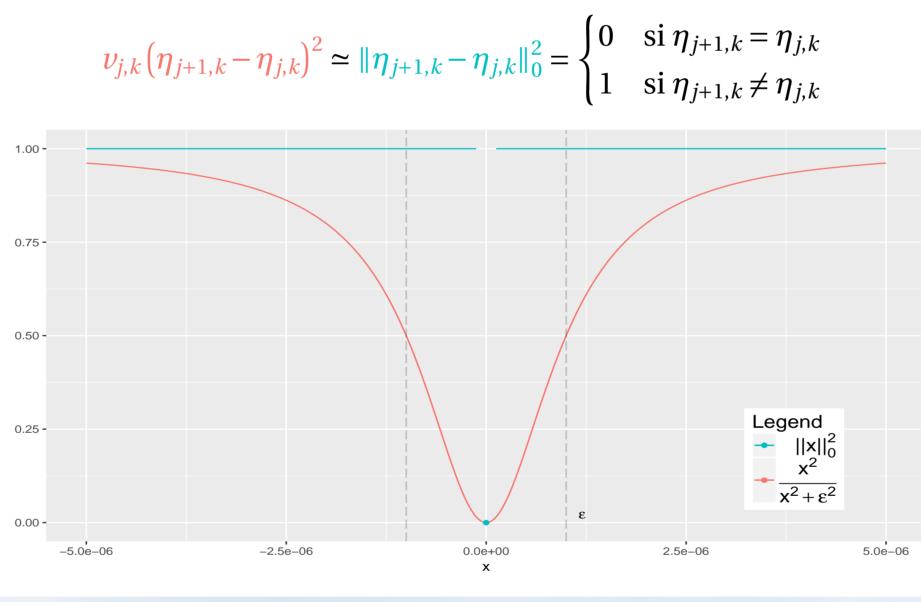
pen is a penalty constant

Types of regularization

- ► Ridge Regularisation L_2 norm with v = w = 1 → Smoothed estimation ► Regularisation L_2 with the iterative **Adaptive Ridge** ^[3] procedure → Segmented estimation
- ► Regularisation L_0 with the iterative Adaptive Ridge ^[3] procedure \rightarrow Segmented estimation The weights are iteratively adapted:

$$\left(\begin{array}{c} v_{j,k} = \left(\left(\eta_{j+1,k} - \eta_{j,k} \right)^2 + \varepsilon^2 \right)^{-1} \\ w_{j,k} = \left(\left(\eta_{j,k} - \eta_{j,k-1} \right)^2 + \varepsilon^2 \right)^{-1} \quad \text{with} \quad \varepsilon \ll 1 \end{array} \right)$$

Approximation of the L₀ norm:



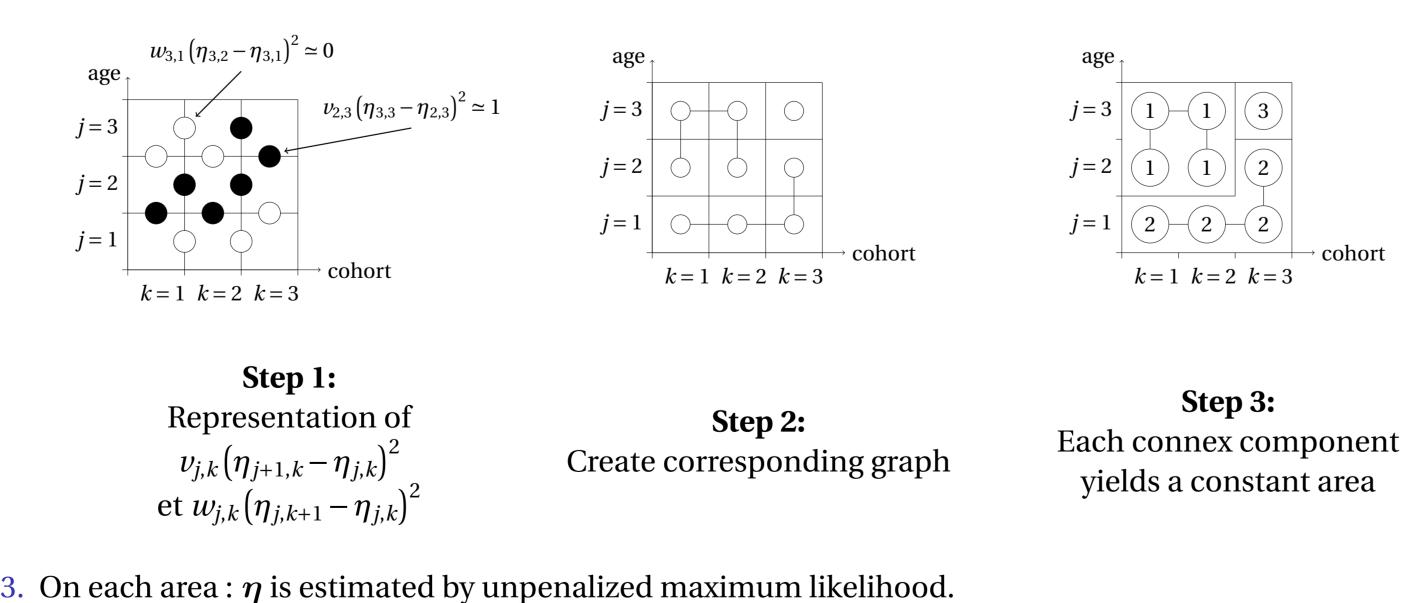
Principle of Model Selection using the L_0 norm

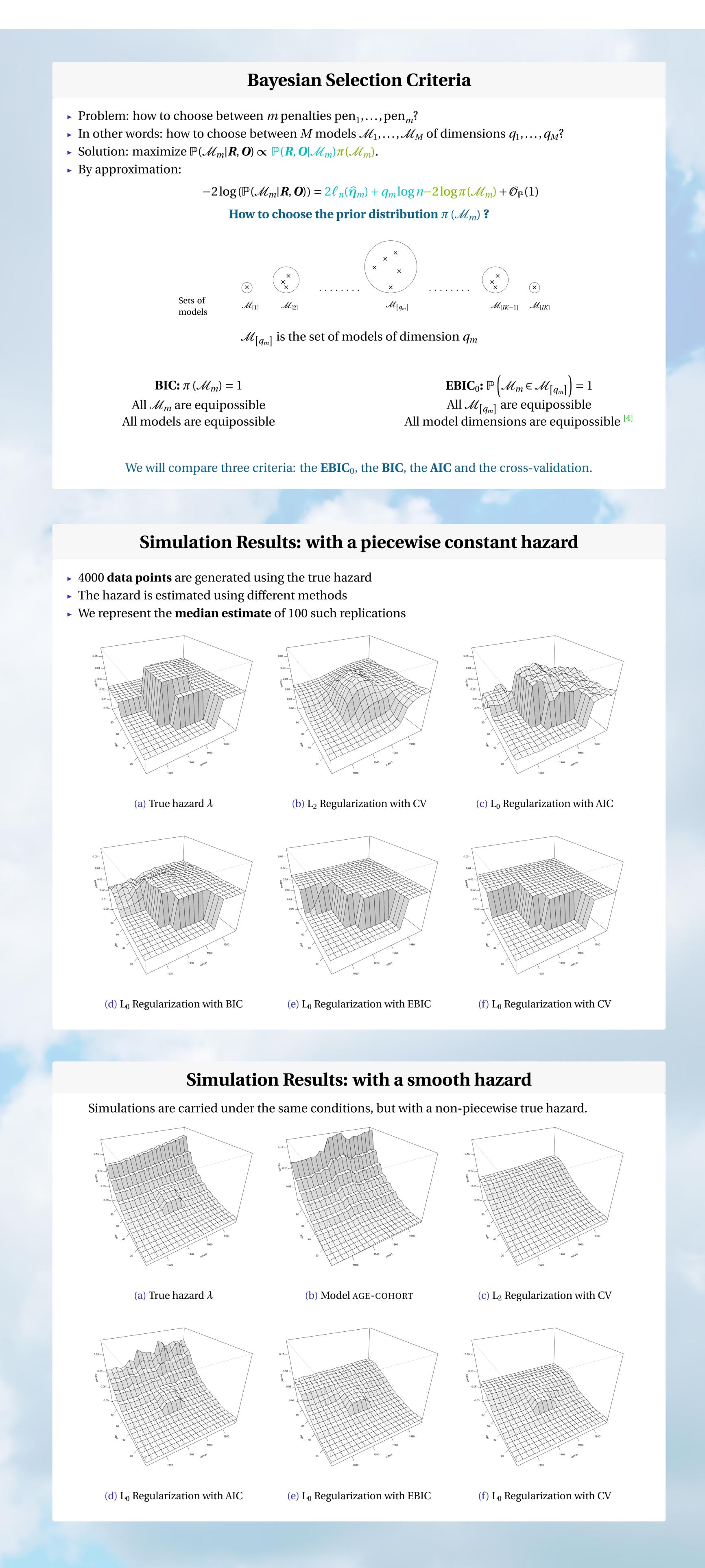
1. We alternate until convergence between:

• Minimize $\ell_n^{\text{pen}}(\eta)$ for fixed **v** and **w**.

• Adapt **v** and **w** using η .

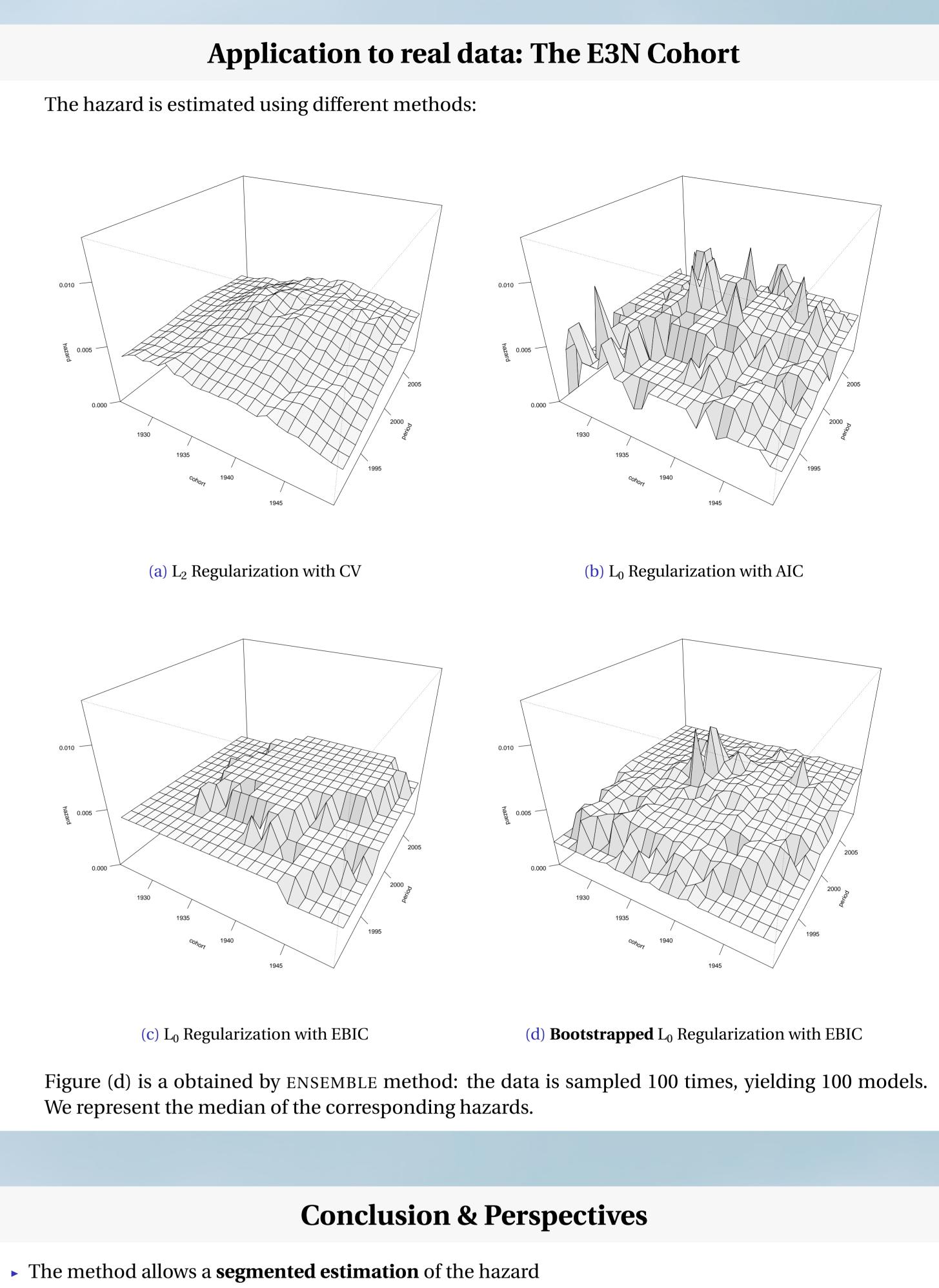
2. Then, the weighted differences of η are used to select areas over which the hazard is constant:











- EBIC₀ more efficient than AIC or BIC
- The model can be extended:

 $\log \lambda_{j,k} = \underbrace{\alpha_j}_{\text{age effet}} + \underbrace{\beta_k}_{\text{cohort effect}} + \underbrace{\delta_{j,k}}_{\text{interaction term}}$

with regularization over $\delta_{j,k}$.

References & Ackownledgment

References

- [1] F. Clavel-Chapelon et al, *Cohort profile: the French E3N cohort study*. International journal of epidemiology, 2014.
- [2] O. Bouaziz and G. Nuel, *L0 Regularization for the Estimation of Piecewise Constant Hazard Rates in Survival Analysis*. Applied Mathematics, 2017.
- [3] F. Frommlet and G. Nuel, An Adaptive Ridge Procedure for L0 Regularization. PloS one, 2016.
- [4] J. Chen and J. Chen, *Extended Bayesian information criteria for model selection with large model spaces*. Biometrika, 20008.

Acknowledgment

This study was realized using data from the Inserm E3N cohort that has been established and is maintained with the financial support of the MGEN, the Gustave Roussy Institute and La Ligue contre le Cancer.