# AN ITERATIVE REGULARIZED METHOD FOR SEGMENTATION WITH APPLICATIONS TO STATISTICS

Vivien Goepp

under the supervision of Pr. O. Bouaziz and Pr. G. Nuel

27 septembre 2019

MAP5, Université Paris Descartes

• Illustrative example : the linear model

$$oldsymbol{y} = oldsymbol{X}oldsymbol{eta}^* + oldsymbol{arepsilon},$$

with  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta^* \in \mathbb{R}^p$ , and  $\mathbb{E}[\varepsilon_i] = 0$ .

• Illustrative example : the linear model

$$y = X\beta^* + \varepsilon,$$

with  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta^* \in \mathbb{R}^p$ , and  $\mathbb{E}[\varepsilon_i] = 0$ .

• Ordinary least square estimate (OLS) :

$$\hat{oldsymbol{eta}}^{\mathsf{ols}} = \mathsf{argmin}_{oldsymbol{eta}} \|oldsymbol{y} - oldsymbol{X}oldsymbol{eta}\|^2 = (oldsymbol{X}^Toldsymbol{X})^{-1}oldsymbol{X}^Toldsymbol{y}$$

• Illustrative example : the linear model

$$y = X\beta^* + \varepsilon,$$

with  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta^* \in \mathbb{R}^p$ , and  $\mathbb{E}[\varepsilon_i] = 0$ .

• Ordinary least square estimate (OLS) :

$$\hat{\boldsymbol{\beta}}^{\mathsf{ols}} = \operatorname{argmin}_{\boldsymbol{\beta}} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \|^2 = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y}$$
  
 $\rightarrow \mathsf{regularize} \ \hat{\boldsymbol{\beta}}^{\mathsf{ols}}$ 

Illustrative example : the linear model

$$y = X\beta^* + \varepsilon,$$

with  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta^* \in \mathbb{R}^p$ , and  $\mathbb{E}[\varepsilon_i] = 0$ .

• Ordinary least square estimate (OLS) :

$$\hat{oldsymbol{eta}}^{\mathsf{ols}} = \mathsf{argmin}_{oldsymbol{eta}} \|oldsymbol{y} - oldsymbol{X}oldsymbol{eta}\|^2 = (oldsymbol{X}^Toldsymbol{X})^{-1}oldsymbol{X}^Toldsymbol{y}$$

ightarrow regularize  $\hat{oldsymbol{eta}}^{\mathsf{ols}}$ 

**Predictive** approach : Ridge regularization :

$$\operatorname{argmin}_{\boldsymbol{\beta}} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \|^2 + \kappa \| \boldsymbol{\beta} \|^2$$

**Explicative** approach : Lasso regularization :

$$\operatorname{argmin}_{\boldsymbol{\beta}} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \|^{2} + \kappa \| \boldsymbol{\beta} \|_{1}$$

1

- The adaptive ridge
- Segmentation in survival analysis
  - Bidimensional hazard rate estimation
  - Extension to age-period-cohort analysis
- Spline regression with knot selection

# 1. The adaptive ridge

#### Penalized likelihood-based variable selection methods

- The LASSO<sup>1</sup> estimate is sparse.
- Definition :

$$\boldsymbol{\beta}^{\text{lasso}} = \text{argmin}_{\boldsymbol{\beta}} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \kappa \|\boldsymbol{\beta}\|_1 \quad (\kappa > 0)$$

<sup>1</sup> :Tibshirani, R., Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society*, 1996.

#### Penalized likelihood-based variable selection methods

- The LASSO<sup>1</sup> estimate is sparse.
- Definition :

$$\boldsymbol{\beta}^{\text{lasso}} = \text{argmin}_{\boldsymbol{\beta}} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \kappa \|\boldsymbol{\beta}\|_1 \quad (\kappa > 0)$$

• Equivalent definition :

$$\min_{\boldsymbol{\beta}} \|\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}\|^2 \quad \text{s.t.} \quad \|\boldsymbol{\beta}\|_1 \leq t \quad (t > 0)$$

<sup>1</sup> :Tibshirani, R., Regression Shrinkage and Selection via the Lasso, *Journal of the Royal Statistical Society*, 1996.

$$\min_{\boldsymbol{\beta}} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 \quad \text{s.t.} \quad \|\boldsymbol{\beta}\|_1 \le t \quad (t > 0)$$

# Illustration with p = 2:



$$\min_{\boldsymbol{\beta}} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 \quad \text{s.t.} \quad \|\boldsymbol{\beta}\|_1 \le t \quad (t > 0)$$

# Illustration with p = 2:



# Bridge penalty : $L_q$ quasi-norm



 $L_q$  norm penalty with q = 2/3

# Bridge penalty : $L_q$ quasi-norm



 $L_q$  norm penalty with q = 2/3

# Bridge penalty : $L_q$ quasi-norm





 $L_q$  norm penalty with q = 2/3

# Approximating $L_q$ norms with the $L_2$ norm

For any  $q \in (0,2)$ , we have<sup>2</sup> :

$$\frac{1}{q} \|\boldsymbol{\beta}\|_{q}^{q} = \inf_{\boldsymbol{w} \in \mathbb{R}_{+}^{d}} \left\{ \frac{1}{2} \sum_{j=1}^{d} w_{j} \beta_{j}^{2} + \frac{2-q}{2q} \sum_{j=1}^{d} |w_{j}|^{\frac{q}{q-2}} \right\}$$

<sup>2</sup> : Mairal, J, Bach, F, and Ponce, J., Sparse Modeling for Image and Vision Processing, Foundations and Trends in Computer Graphics and Vision, 2014.

#### Approximating $L_q$ norms with the $L_2$ norm

For any  $q \in (0,2)$ , we have<sup>2</sup> :

$$\frac{1}{q} \|\boldsymbol{\beta}\|_{q}^{q} = \inf_{\boldsymbol{w} \in \mathbb{R}_{+}^{d}} \left\{ \frac{1}{2} \sum_{j=1}^{d} w_{j} \beta_{j}^{2} + \frac{2-q}{2q} \sum_{j=1}^{d} |w_{j}|^{\frac{q}{q-2}} \right\}$$



<sup>2</sup> : Mairal, J, Bach, F, and Ponce, J., Sparse Modeling for Image and Vision Processing, Foundations and Trends in Computer Graphics and Vision, 2014.

The problem becomes

$$\min_{\boldsymbol{\beta}} \left\{ \|\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}\|^2 + \kappa \|\boldsymbol{\beta}\|_q^q \right\} = \min_{\boldsymbol{\beta}} \inf_{\boldsymbol{w} \in \mathbb{R}^d_+} \left\{ \ell(\boldsymbol{\beta}, \boldsymbol{w}) \right\}$$

with

$$\ell(\boldsymbol{\beta}, \boldsymbol{w}) = \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \frac{\kappa q}{2} \sum_j w_j \beta_j^2 + \frac{\kappa(2-q)}{2} \sum_j |w_j|^{\frac{q}{q-2}}$$

The L<sub>q</sub> adaptive ridge : **minimize alternatively** 

• 
$$\operatorname{arg\,min}_{\boldsymbol{w}} \left\{ \ell(\boldsymbol{\beta}, \boldsymbol{w}) \right\} = (\left|\beta_j\right|^2)^{\frac{q-2}{2}}$$

• 
$$\arg\min_{\boldsymbol{\beta}} \left\{ \ell(\boldsymbol{\beta}, \boldsymbol{w}) \right\} = \arg\min_{\boldsymbol{\beta}} \left\| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \right\|^2 + \frac{\kappa q}{2} \sum_j w_j \beta_j^2$$

The problem becomes

$$\min_{\boldsymbol{\beta}} \left\{ \|\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}\|^2 + \kappa \|\boldsymbol{\beta}\|_q^q \right\} = \min_{\boldsymbol{\beta}} \inf_{\boldsymbol{w} \in \mathbb{R}^d_+} \left\{ \ell(\boldsymbol{\beta}, \boldsymbol{w}) \right\}$$

with

$$\ell(\boldsymbol{\beta}, \boldsymbol{w}) = \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \frac{\kappa q}{2} \sum_j w_j \beta_j^2 + \frac{\kappa(2-q)}{2} \sum_j |w_j|^{\frac{q}{q-2}}$$

The L<sub>q</sub> adaptive ridge : **minimize alternatively** 

• 
$$\operatorname{arg\,min}_{\boldsymbol{w}} \left\{ \ell(\boldsymbol{\beta}, \boldsymbol{w}) \right\} = \left( |\beta_j|^2 + \delta^2 \right)^{\frac{q-2}{2}}$$

• 
$$\arg\min_{\boldsymbol{\beta}} \left\{ \ell(\boldsymbol{\beta}, \boldsymbol{w}) \right\} = \arg\min_{\boldsymbol{\beta}} \left\| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta} \right\|^2 + \frac{\kappa q}{2} \sum_j w_j \beta_j^2$$

- The previous equation is defined when q = 0!
- The  $\mathsf{L}_0$  adaptive ridge estimate is the limit of the iterations :
  - $w_j \leftarrow (\beta_j^2 + \delta^2)^{-1}$
  - $\boldsymbol{\flat} \quad \boldsymbol{\beta} \leftarrow \operatorname{argmin}_{\boldsymbol{\beta}} \| \boldsymbol{y} \boldsymbol{X} \boldsymbol{\beta} \|^2 + \kappa \sum_j w_j \beta_j^2$
- But it does not correspond to the L<sub>0</sub> penalty.

# The algorithm of the L<sub>0</sub> adaptive ridge

· It corresponds to the minimization of

$$\min_{\boldsymbol{\beta}} \|\boldsymbol{y} - \boldsymbol{X} \boldsymbol{\beta}\|^2 + \kappa \sum_j \log(\beta_j^2 + \delta^2)$$

using the MM optimization :



- Define  $q(\beta_j|\beta_j^{(k)}) = \log(\beta_j^{(k)2} + \delta^2) + (\beta_j^2 \beta_j^{(k)2})$
- Iterate over  $\beta^{(k+1)} \leftarrow \operatorname{argmin}_{\beta} q(\beta|\beta^{(k)})$ .

Define  $\ell(\beta) = -\log(L(\beta))$ . L<sub>0</sub> adaptive ridge :

$$\begin{split} & \boldsymbol{w}^{(0)} \leftarrow \mathbf{1} \\ & \boldsymbol{k} \leftarrow 1 \\ & \mathbf{do} \\ & \boldsymbol{\beta}^{(k)} \leftarrow \mathrm{argmin}_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}) + \kappa \sum_{j} w_{j}^{(k-1)} \beta_{j}^{(k)2} \\ & w_{j}^{(k)} \leftarrow (\beta_{j}^{(k)2} + \delta^{2})^{-1} \\ & \boldsymbol{k} \leftarrow k + 1 \\ & \mathbf{while} \ |w_{j}^{(k)} \beta_{j}^{(k)2} - 0.5| < 0.5 - \varepsilon \\ & \text{For every } j, \text{ we set } w_{j}^{(k)} \beta_{j}^{(k)2} \in \{0, 1\} \\ & \text{ If } w_{j}^{(k)} \beta_{j}^{(k)2} = 1 \text{, we set } \hat{\boldsymbol{\beta}}_{j} \leftarrow \hat{\boldsymbol{\beta}}_{j}^{\text{mle}} \end{split}$$

Fused  $L_0$  adaptive ridge :

- Segmentation of  $\beta$ .
- Penalize over  $(\beta_j \beta_{j-1})_j$ .

Fused L<sub>0</sub> adaptive ridge :

- Segmentation of  $\beta$ .
- Penalize over  $(\beta_j \beta_{j-1})_j$ .

$$\begin{split} & \boldsymbol{w}^{(0)} \leftarrow \mathbf{1} \\ & k \leftarrow 1 \\ & \mathbf{do} \\ & \boldsymbol{\beta}^{(k)} \leftarrow \operatorname{argmin}_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}) + \kappa \sum_{j} w_{j}^{(k-1)} (\boldsymbol{\beta}_{j}^{(k)} - \boldsymbol{\beta}_{j-1}^{(k)})^{2} \\ & w_{j}^{(k)} \leftarrow ((\boldsymbol{\beta}_{j}^{(k)} - \boldsymbol{\beta}_{j-1}^{(k)})^{2} + \delta^{2})^{-1} \\ & k \leftarrow k + 1 \\ & \text{while } |w_{j}^{(k)} (\boldsymbol{\beta}_{j}^{(k)} - \boldsymbol{\beta}_{j-1}^{(k)})^{2} - 0.5| < 0.5 - \varepsilon \end{split}$$

For every j, we set  $w_j^{(k)}(\beta_j^{(k)} - \beta_{j-1}^{(k)})^2 \in \{0, 1\}$ On each segment :  $\hat{\beta} \leftarrow \hat{\beta}^{mle}$ 

#### Comparison : smoothing vs fused L<sub>0</sub> adaptive ridge



Fused ridge : Each  $\kappa$  yields an *estimate* 



 $L_0$  adaptive ridge segmentation : Each  $\kappa$  yields a *model*  Selected papers :

- Grandvalet (1998) : Approximate L<sub>1</sub> norm
- Daubechies et al (2010) :
  - ► Study of the case *q* > 0 in sparse sensing
  - Speed of convergence is superlinear for q > 0.
- Frommlet and Nuel (2016) : Introduced L $_q$  ( $q \ge 0$ ) Adaptive Ridge
- Liu et al (2017) : The L\_0 adaptive ridge : empirical study in high dimension ( $p \gg n$ ).
- Dai et al (2018) : The L<sub>0</sub> adaptive ridge regression has the oracle properties.

2. Application to survival analysis

2.1 Bidimensional estimation of the hazard rate

# Motivating application : the SEER data

#### Study of mortality following breast cancer diagnosis

- US registry dataset of breast cancer
- Primary, unilateral, malignant and invasive cancers
- 1.2 million of patients, 60% of censoring
- The dates of diagnoses range from 1973 to 2014
- Variable of interest : the time from diagnosis until death (from cancer).





- Variable of interest :  $T^*$ , time before an event of interest.
- But we don't observe the  $T_i^*$ s, but

$$\begin{cases} T_i = \min\left(T_i^*, C_i\right) \\ \Delta_i = \mathbb{1}_{T_i = T_i^*}. \end{cases}$$

#### **Right-censored data**



- Variable of interest : T\*, time before an event of interest.
- But we don't observe the  $T_i^*$ s, but

$$\begin{cases} T_i = \min\left(T_i^*, C_i\right), \\ \Delta_i = \mathbb{1}_{T_i = T_i^*}. \end{cases}$$

.

- $T_i$ : observed time  $C_i$ : censoring
- $U_i$ : covariate
- Aim : infer the hazard rate :

$$\lambda(t|u) = \lim_{\delta t \to 0} \frac{\mathbb{P}\left(t \le T^* \le t + \delta t | T^* > t, U = u\right)}{\delta t}$$

• **Our approach** : discretize  $\lambda$  :

$$\lambda(t|u) = \sum_{j=1}^{J} \sum_{k=1}^{K} \lambda_{j,k} I_{[c_{j-1},c_j) \times [d_{k-1},d_k)}(t,u)$$

with segmentation over  $\lambda_{j,k}$ .

Exhaustive statistics :

- $O_{j,k}$  : number of observed events in rectangle (j,k)
- $R_{j,k}$  : total time at risk in rectangle (j,k)The negative log-likelihood writes :

$$\ell_n(\eta) = \sum_{j=1}^J \sum_{k=1}^K \{ R_{j,k} \exp(\eta_{j,k}) - \eta_{j,k} O_{j,k} \} \text{ with } \log \lambda_{j,k} = \eta_{j,k}.$$

MLE is explicit :  $\lambda_{j,k} = O_{j,k}/R_{j,k}$ .

.

Exhaustive statistics :

•  $O_{j,k}$  : number of observed events in rectangle (j,k)

•  $R_{j,k}$  : total time at risk in rectangle (j,k)The negative log-likelihood writes :

$$\ell_n(\eta) = \sum_{j=1}^J \sum_{k=1}^K \{ R_{j,k} \exp(\eta_{j,k}) - \eta_{j,k} O_{j,k} \} \text{ with } \log \lambda_{j,k} = \eta_{j,k}.$$

MLE is explicit :  $\lambda_{j,k} = O_{j,k}/R_{j,k}$ . Bidimensional fused adaptive ridge :

$$\ell_{n}^{\mathsf{pen}}(\boldsymbol{\eta}) = \ell_{n}(\boldsymbol{\eta}) + \kappa \sum_{j,k} \left\{ v_{j,k} \left( \eta_{j+1,k} - \eta_{j,k} \right)^{2} + w_{j,k} \left( \eta_{j,k+1} - \eta_{j,k} \right)^{2} \right\}.$$

#### Segmentation into constant areas







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

(b) Corresponding graph

(c) Segmentation into connected components

#### Segmented estimate of the hazard rate



#### (Unpenalized) MLE

#### Segmented estimate of the hazard rate





#### (Unpenalized) MLE

L<sub>2</sub> Regularization

### **Results with segmentation**







Fused Adaptive Ridge (greyscale)
# 2.2 Extension to age-period-cohort analysis

## Age-period-cohort analysis



# Age-period-cohort analysis



Key relation : period = age + cohort.

We want to infer the effects of age, period, and cohort.

Example : incidence of breast cancer :

- age effect : menopause
- cohort effect : carcinogenic baby food
- period effect : nuclear accident

We define one parameter vector for each effect :  $\alpha$ ,  $\beta$  et  $\gamma$ 

1. In the AGE-PERIOD-COHORT  $^3$  model, we assume

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

- Non-identifiable :
  - infer  $\Delta^2 \alpha$ ,  $\Delta^2 \beta$  et  $\Delta^2 \gamma$ .
  - or add a constraint to the model.
- 2. In the  $\ensuremath{\mathsf{AGE}}\xspace{-}\ensuremath{\mathsf{COHORT}}\xspace^3$  model, we assume

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

- Outer product structure  $\rightarrow$  regularizing
- Additive effect of the variables : strong a priori

 $^3$  : Carstensen, B., Age-period-cohort models for the Lexis diagram, Statistics in medicine, 2007.

• We introduce an AGE-COHORT-INTERACTION model :

$$\log\left(\lambda_{j,k}\right) = \alpha_j + \beta_k + \delta_{j,k},$$

where  $\delta_{j,k}$  is the interaction (with  $\delta_{1,k} = \delta_{j,1} = 0$ ).

• We regularize over the differences of  $\delta_{j,k}$  :

$$\ell_n^{\mathsf{pen}}(\boldsymbol{\theta}) = \ell_n(\boldsymbol{\theta}) + \kappa \sum_{j,k} \left\{ v_{j,k} \left( \delta_{j+1,k} - \delta_{j,k} \right)^2 + w_{j,k} \left( \delta_{j,k+1} - \delta_{j,k} \right)^2 \right\}.$$

- Compromise between :
  - AC Model :  $\kappa \to \infty$
  - MLE :  $\kappa \to 0$

- J = 20 age intervals and K = 20 cohort intervals
- Sample the cohort uniformly
- Sample the age using the hazard rate  $(\lambda_{j,k})$
- Uniform censoring over the age [75, 100]
- Infer  $(\alpha, \beta, \delta)$  in the ACI model.
- We represent medians over 100 repetitions





True hazard  $\lambda_{j,k}^*$ 

Age-cohort model ( $\log \lambda_{j,k} = \alpha_j + \beta_k$ )



True hazard  $\lambda_{j,k}^*$ 

ACI model : estimated  $\hat{\lambda}_{j,k}$ 

## **Results with the ACI model**





True interaction :  $\delta_{j,k}^*$ 

ACI model : estimated  $\hat{\delta}_{j,k}$ 

## **Results with the ACI model**



# 3. Spline Regression with Automatic Knot Selection

Let  $(x_i, y_i) \in \mathbb{R}$  :

$$y_i = f(x_i) + \varepsilon_i, \quad 1 \le i \le n,$$

- f "smooth"
- $\varepsilon_i$  i.i.d and centered
- $x_i \in [a, b]$

Aim : infer f.

- Define the spline order  $q \ge 1$  and knots  $(t_1, \cdots, t_k) \in [a, b]$ .
- The estimate is the spline  $\hat{f}(x) = \sum_{j=1}^{q+k} a_j B_{j,q}(x)$ , where  $B_{j,q}$  is the B-spline of order q.
- We minimize in  $oldsymbol{a} \in \mathbb{R}^{q+k}$  :

SS
$$(\boldsymbol{a}, \boldsymbol{t}) = \sum_{i=1}^{n} \left\{ y_i - \sum_{j=1}^{q+k} a_j B_{j,q}(x_i) \right\}^2 = \|\boldsymbol{y} - \boldsymbol{B}\boldsymbol{a}\|^2,$$

# **B-spline regression**

 $(B_{j,q})_{j=1}^{k+q}$  is a basis of splines of order q.



B-splines bases with three knots : (0.25, 0.5, 0.75).

#### **Knot placement**

#### Where to place knots?



Uniform knots

Helmet crash test data<sup>5</sup>

<sup>5</sup> : Silverman, B.W., Some Aspects of the Spline Smoothing Approach to Non-Parametric Regression Curve Fitting, *Journal of the Royal Statistical Society*, 1985.

#### **Knot placement**

#### Where to place knots?



Helmet crash test data<sup>5</sup>

<sup>5</sup> : Silverman, B.W., Some Aspects of the Spline Smoothing Approach to Non-Parametric Regression Curve Fitting, *Journal of the Royal Statistical Society*, 1985.





Uniform = quantile knots



#### It is only based on the *x*-distribution of the points :

- Need to find optimal knot position
- Their optimal position is informative about breakpoints in f.

#### Two different approaches

- **Predictive approach** : place too many knots and regularize *a* : P-splines<sup>4</sup>.
- Explicative approach :
  - ► Jointly optimize w.r.t. t and a : Jupp (1978), Lindstrom (1999)
  - Monte-Carlo-based Approach : Denison et al. (1998), DiMatteo et al. (2001)

<sup>4</sup> : Eilers, P. and Marx, D., Flexible Smoothing with B-splines and Penalties, *Statistical Science*, 1996.

## **A-Spline**

Our approach :

- Set many initial knots
- · Successively remove the least relevant knots
- Using the fused adaptive ridge

Enforce the successive values of a to be equal.

$$\begin{split} \operatorname{PSS}\left(\boldsymbol{a},\boldsymbol{w}\right) &= \|\boldsymbol{y} - \boldsymbol{B}\boldsymbol{a}\|_{2}^{2} + \kappa \sum_{j=q+1}^{q+k} w_{j} \left(\Delta^{q} a_{j}\right)^{2}, \\ \text{with} \quad \Delta a_{j} &= a_{j} - a_{j-1} \quad \text{and} \quad \Delta^{q} = \Delta^{q-1} \circ \Delta \end{split}$$

## **A-Spline**

Our approach :

- Set many initial knots
- · Successively remove the least relevant knots
- Using the fused adaptive ridge

Enforce the successive values of a to be equal.

$$PSS(\boldsymbol{a}, \boldsymbol{w}) = \|\boldsymbol{y} - \boldsymbol{B}\boldsymbol{a}\|_{2}^{2} + \kappa \sum_{j=q+1}^{q+k} w_{j} (\Delta^{q} a_{j})^{2},$$
  
with  $\Delta a_{j} = a_{j} - a_{j-1}$  and  $\Delta^{q} = \Delta^{q-1} \circ \Delta$ 

Explicit iteration step :

$$\begin{aligned} \operatorname{argmin}_{\boldsymbol{a}} \operatorname{PSS}(\boldsymbol{a}, \boldsymbol{w}) &= (\boldsymbol{B}^T \boldsymbol{B} + \kappa \boldsymbol{D}_{\boldsymbol{q}}^T \boldsymbol{W} \boldsymbol{D}_{\boldsymbol{q}})^{-1} \boldsymbol{B}^T \boldsymbol{y} \\ & \text{with} \quad \boldsymbol{D}_1 = \begin{bmatrix} 1 & -1 & & \\ & 1 & -1 & \\ & \ddots & \ddots & \\ & & & 1 & -1 \end{bmatrix} \end{aligned}$$

Illustration on simulated data :

- $f(x) = \sin(6\pi x) * 0.5 + 0.5$
- $\varepsilon \sim \mathcal{N}(0, 0.15^2)$  and n = 200.
- q = 4
- 40 uniform initial knots





A-splines

## Results on real data : detecting changes in mean



aCGH profile of bladder tumor samples<sup>6</sup> :

<sup>6</sup> : Stransky, N. et al, Regional Copy Number–Independent Deregulation of Transcription in Cancer, *Nature Genetics*, 2006.

<sup>7</sup> : Killick, R. et al, Optimal detection of changepoints with a linear computational cost, *Journal of the American Statistical Association*, 2006.

## Results on real data : detecting changes in mean



aCGH profile of bladder tumor samples<sup>6</sup> :

<sup>6</sup> : Stransky, N. et al, Regional Copy Number–Independent Deregulation of Transcription in Cancer, *Nature Genetics*, 2006.

<sup>7</sup> : Killick, R. et al, Optimal detection of changepoints with a linear computational cost, *Journal of the American Statistical Association*, 2006.

## Results on real data : detecting changes in mean



aCGH profile of bladder tumor samples<sup>6</sup> :

<sup>6</sup> : Stransky, N. et al, Regional Copy Number–Independent Deregulation of Transcription in Cancer, *Nature Genetics*, 2006.

<sup>7</sup> : Killick, R. et al, Optimal detection of changepoints with a linear computational cost, *Journal of the American Statistical Association*, 2006.

## Results on real data : detecting changes of slope

Light measurement data set<sup>8</sup> :



<sup>8</sup> : Sigrist, M. et al, Air monitoring by spectroscopic techniques, *John Wiley & Sons*, 1994.

<sup>9</sup> : Friedman ; J., Multivariate Adaptive Regression Splines, *Journal of the American Statistical Association*, 1991.

## Results on real data : detecting changes of slope

Light measurement data set<sup>8</sup> :



<sup>8</sup> : Sigrist, M. et al, Air monitoring by spectroscopic techniques, *John Wiley & Sons*, 1994.

<sup>9</sup> : Friedman ; J., Multivariate Adaptive Regression Splines, *Journal of the American Statistical Association*, 1991.

## Results on real data : detecting changes of slope

Light measurement data set<sup>8</sup> :



<sup>8</sup> : Sigrist, M. et al, Air monitoring by spectroscopic techniques, *John Wiley & Sons*, 1994.

<sup>9</sup> : Friedman ; J., Multivariate Adaptive Regression Splines, *Journal of the American Statistical Association*, 1991.

#### Results on real data : splines of higher order

Number of disasters in coal mines in the UK<sup>10</sup> :



<sup>10</sup>: Diggle, P and Marron, J., Equivalence of Smoothing Parameter Selectors in Density and Intensity Estimation, JASA, 1988.

#### Results on real data : splines of higher order

Number of disasters in coal mines in the  $UK^{10}$ :



<sup>10</sup>: Diggle, P and Marron, J., Equivalence of Smoothing Parameter Selectors in Density and Intensity Estimation, JASA, 1988.

#### Results on real data : splines of higher order

Number of disasters in coal mines in the UK<sup>10</sup> :



Poisson model :  $\mathbb{E}[\boldsymbol{y}|\boldsymbol{x}] = \boldsymbol{\mu}$  with  $\hat{\boldsymbol{\mu}} = \exp(\boldsymbol{B}\hat{\boldsymbol{a}})$ 

<sup>10</sup> : Diggle, P and Marron, J., Equivalence of Smoothing Parameter Selectors in Density and Intensity Estimation, JASA, 1988.

# Comparison of mean squarred error (MSE) on simulated data :



## Comparison of mean squarred error (MSE) on simulated data :



45

# Conclusion :

- Papers :
  - Bidimensional estimation of the hazard rate [submitted]
  - Spline regression with automatic knot selection [submitted]
  - Age-cohort-interaction model [under preparation]
  - Segmentation of geographic-based data [under preparation]
- Communications :
  - Conference talks : SAM 2017, IWAP 2018
  - Conference posters : SMPGD 2018, IBC 2018, SMPGD 2019
  - Three invited seminars
- Three R packages : hazreg, aspline, and graphseg.

Perspectives :

- Study of consistency of the segmentation (ACI model).
- Use of splines in age-period-cohort models.
- Application to study of evolution of mortality causes in France.
# Annex 1 : Adaptive ridge and MM optimization

Local Quadratic Approximation : Solve iteratively

$$\min_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}) + \kappa \sum_{j} p(|\beta_{j}|),$$

where p is even, concave on  $\mathbb{R}_+$ , and nondecreasing.

$$p(|\beta_j^{(k)}|) + (\beta_j^2 - \beta_j^{(k)2}) \frac{p'(|\beta_j^{(k)}|)}{|\beta_j^{(k)}|} \succeq p(|\beta_j|)$$

### MM Optimization :

- $\min_{\beta} \ell^{\mathsf{pen}}(\beta|\beta^{(k)})$
- Update  $oldsymbol{eta}^{(k)}$  to  $oldsymbol{eta}^{(k+1)}$

The sequence  $\beta^{(k)}$  tends to a local minimum of  $\ell^{\text{pen}}(\beta)$ .

Fan, J. and Li, R., Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties, *Journal of the American Statistical Association*, 1996.

When  $p(|\beta_j|) = \log(\beta_j^2 + \delta^2)$ , this is the L<sub>0</sub> adaptive ridge.



#### The Adaptive Lasso comes from the LLA of a penalty.

$$p(|\beta_j|) + (|\beta_j| - |\beta_j^{(k)}|)p'(|\beta_j^{(k)}|) \succeq p(|\beta_j|)$$

## Annex 2 : Connexion with adaptive lasso (2)

Example : with  $p(|\beta_j|) = \log(|\beta_j|)$ .



- Each iteration is sparse
- Zou and Li (2008) offer to stop at one iteration

Compared with the adaptive ridge :

- It takes less iterations to converge.
- Each iteration is slower : gradient methods versus newton methods.

- Problem : choose between M models  $\mathcal{M}_1, \ldots, \mathcal{M}_M$  of dimensions  $q_1, \ldots, q_M$ .
- Solution : maximize  $\mathbb{P}(\mathcal{M}_m | \mathbf{R}, \mathbf{O}) \propto \mathbb{P}(\mathbf{R}, \mathbf{O} | \mathcal{M}_m) \pi(\mathcal{M}_m)$ .
- By approximation :

 $-2\log\left(\mathbb{P}(\mathcal{M}_m|\mathbf{R}, \mathbf{O})\right) = 2\ell_n(\widehat{\boldsymbol{\eta}}_m) + q_m\log n - 2\log \pi(\mathcal{M}_m) + \mathcal{O}_{\mathbb{P}}(1)$ 

• We must choose the prior *a priori*  $\pi(\mathcal{M}_m)$ 

BIC :  $\pi(\mathcal{M}_m) = 1$ All the  $\mathcal{M}_m$  are equiprobable  $\mathsf{EBIC}_0: \mathbb{P}\left(\mathcal{M}_m \in \mathcal{M}_{[q_m]}\right) = 1$ All the  $\mathcal{M}_{[q_m]}$  are equiprobable



 $\mathcal{M}_{[q_m]}$  is the set of models with  $q_m$  parameters

#### We compare different model selection criteria :

1. 
$$\mathsf{BIC}(m) = 2\ell_n(\widehat{\boldsymbol{\eta}}_m) + q_m\log n$$

2. 
$$\mathsf{EBIC}_0(m) = 2\ell_n(\widehat{\eta}_m) + q_m \log n + 2\log \binom{JK}{q_m}$$

3. AIC
$$(m) = 2\ell_n(\widehat{\boldsymbol{\eta}}_m) + 2q_m$$

4. K-fold Cross validation (CV)