# **Binacox:** automatic cut-points detection in high-dimensional Cox model

<sup>1</sup>LPSM, UMR 8001, Sorbonne University, Paris, France <sup>2</sup>Modal'X, UPL, Univ Paris Nanterre, F92000 Nanterre France <sup>3</sup>LaMME, UEVE and UMR 8071, Paris Saclay University, Evry, France <sup>4</sup>Biomedical Informatics and Public Health Department European Georges Pompidou Hospital, Assistance Publique-Hôpitaux de Paris and INSERM UMRS 1138, Centre de Recherche des Cordeliers, Paris, France

### Objectives

- 1 Introduce the cut-points detection problem
- **2**Present the estimation procedure
- **3** Give theoretical guarantees
- 4 Illustrate the method on simulated and real
- high-dimensional data

#### Introduction

Translating significant continuous prognostic biomarkers into clinical decision often requires determining cut-points. We introduce a prognostic method called *Binacox* to deal with this problem of detecting multiple cut-points per features in a multivariate high-dimentional survival setting.

Let us denote T and C the times of the event of interest and censoring times respectively,  $X \in \mathbb{R}^p$ the vector of features,  $Z = T \wedge C$  the right-censored time and  $\Delta = \mathbb{1}(\{T \leq C\})$  the censoring indicator. Assume that intensity of events for patient i is given by

$$\lambda^{\star}(t|X_i = x_i) = \lambda_0^{\star}(t)e^{f^{\star}(x_i)}, \qquad (1)$$

where  $\lambda_0^{\star}(t)$  is the baseline hazard function, and

$$f^{\star}(x_i) = \sum_{j=1}^{p} \sum_{k=1}^{K_j^{\star}+1} \beta_{j,k}^{\star} \mathbb{1}(x_{i,j} \in I_{j,k}^{\star}), \qquad (2)$$

with  $I_{j,k}^{\star} = (\mu_{j,k-1}^{\star}, \mu_{j,k}^{\star}]$  for  $k \in \{1, \dots, K_{j}^{\star} + 1\}$ . Our goal is to estimate simultaneously  $\mu^*$  and  $\beta^*$ .



Fig.1: Illustration of data simulated according to (1).

# Simon Bussy<sup>1</sup>, Mokhtar Z. Alaya<sup>2</sup>, Agathe Guilloux<sup>3</sup> and Anne-Sophie Jannot<sup>4</sup>

#### Binarization

First, we one-hot encode all features to obtain  

$$\begin{aligned} x_i^B &= (x_{i,1,1}^B, \dots, x_{i,1,d_1+1}^B, \dots, x_{i,p,d_p+1}^B)^\top, \\ \text{with } x_{i,j,l}^B &= \begin{cases} 1 & \text{if } x_{i,j} \in I_{j,l}, \\ 0 & \text{otherwise,} \end{cases} \text{ and where } I_{j,l} &= \\ \begin{pmatrix} 1 & \text{if } x_{i,j} \in I_{j,l}, \\ 0 & \text{otherwise,} \end{cases} \text{ and where } I_{j,l} &= \\ \begin{pmatrix} \mu_{j,l-1}, \mu_{j,l} \end{bmatrix} \text{ with } \mu_{j,l} &= l/(d_j+1) \text{ for instance.} \\ \text{Hence, we define} \end{cases} \text{ bina}(\beta) &= \sum_{j=1}^p \left( \sum_{l=2}^{d_j+1} \hat{\omega}_{j,l} |\beta_{j,l} - \beta_{j,l-1}| + \delta_1(\beta_{j,\bullet}) \right), \\ \text{with} &= \begin{cases} 0 & \text{if } \mathbf{1}^\top u = 0, \\ \infty & \text{otherwise.} \end{cases} \text{ with} \end{cases} \text{ and otherwise,} &= \sum_{j=1}^{p} f_{\beta_{j,\bullet}}(x_i) \\ \text{where } f_{\beta_{j,\bullet}}(x_i) &= \sum_{j=1}^{d_j+1} \beta_{j,l} \mathbbm{1}(x_{i,j} \in I_{j,l}). \text{ Thus, } f_{\beta} \text{ is a candidate for the ortimation of } f^* \text{ in } (2) \end{cases}$$

a candidate for the estimation of  $J^{\uparrow}$  in (2). We obtain the binarized partial negative loglikelihood

$$\ell_n(f_\beta) = -\frac{1}{n} \sum_{i=1}^n \delta_i \Big\{ f_\beta(x_i) - \log \sum_{i': z_{i'} \ge z_i} e^{f_\beta(x_{i'})} \Big\}.$$

## Fast oracle inequality in prediction

One has 
$$KL_n(f^*, f_{\hat{\beta}}) \leq (1+c_1) \inf_{\substack{\beta \in \mathscr{B}_{p+d} \\ \forall j, \mathbf{1}^\top \beta_{j,\bullet} = 0}} \left\{ KL_n(f^*, f_{\beta}) + c_2 |\mathcal{A}(\beta)| \max_{j=1,\dots,p} \|(\hat{\omega}_{j,\bullet})_{\mathcal{A}_j(\beta)}\|_{\infty}^2 \right\}$$
 with high-probability, where  $\hat{\omega}_{j,l} = \mathcal{O}\left(d_j \sqrt{\pi_n \log(p+d)/n}\right)$  are data-driven weights,  $\pi_n = |\{i=1,\dots,n:\delta_i=1\}$ 

1 | /n, and  $c_1, c_2$  are positive constants (with  $c_2$  resulting from compatibily conditions).

#### Simulation study





 $\mathcal{S}'$  the indexes of features with at least one true cut-point.

#### Estimation procedure

Let  $\mathcal{A}_j(\hat{\beta}) = \{l : \hat{\beta}_{j,l} \neq \hat{\beta}_{j,l-1}\} = \{\hat{l}_{j,1}, \dots, \hat{l}_{j,s_j}\},\$ then one get

$$\widehat{\mu}_{j,\bullet} = (\mu_{j,\hat{l}_{j,1}},\ldots,\mu_{j,\hat{l}_{j,s_j}})^{\top}$$

$$A_{j,i}(\widehat{A}) = \widehat{K}_{i,i}$$

cut-point.

with  $s_j = |\mathcal{A}_j(\beta)| = K_j$ .





Fig.5: Illustration on the top-4 genes of GBM cancer. For instance, the first gene SOD3 is related to an antioxidant enzyme that may protect in particular the brain from oxidative stress, which is believed to play a key role in tumor formation.

Table: Comparison of risk prediction in terms of C-index on three TCGA datasets. Taking into account the detected cut-points significantly improves predictions.

Cancer GBM KIRC BRCA

Our method provides a new way to model nonlinear features associations, and powerful interpretation aspects that could be useful in both clinical research and daily practice: in addition to its raw feature selection ability, the estimated cut-points could directly be used in clinical routine. Software: github.com/SimonBussy/binacox.

Simon Bussy, PhD Student Email: simon.bussy@gmail.com Website: www.simonbussy.com





# **Results on TCGA data**

er	Continuous	Binacox	MT-B	MT-LS	CoxBoost	RSF
I	0.660	0.806	0.753	0.768	0.684	0.691
2	0.682	0.727	0.663	0.663	0.679	0.686
4	0.713	0.849	0.741	0.738	0.723	0.746

# Conclusion

#### **Contact Information**