

## Long Term Survival Probabilities and Kaplan-Meier Estimator

Ion Grama\*

Jean Marie Tricot<sup>†</sup>Jean François Petiot<sup>‡</sup>

### Abstract

The Kaplan-Meier nonparametric estimator has become a standard tool for estimating a survival time distribution in a right censoring schema. However, if the censoring rate is high, this estimator do not provide a reliable estimation of the extreme survival probabilities. In this paper we propose to combine the nonparametric Kaplan-Meier estimator and a parametric-based model into one construction. The idea is to fit the tail of the survival function with a parametric model while for the remaining to use the Kaplan-Meier estimator. A procedure for the automatic choice of the location of the tail based on a goodness-of-fit test is proposed. This technique allows us to improve the estimation of the survival probabilities in the mid and long term. We perform numerical simulations which confirm the advantage of the proposed method.

**Key Words:** Adaptive estimation, censored data, model selection, prediction, survival analysis, survival probabilities

### 1. Introduction

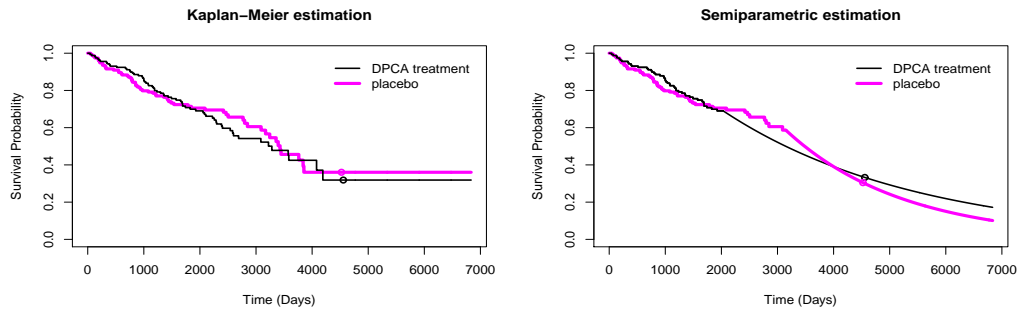
Let  $(X_i, C_i, Z_i)'$ ,  $i = 1, \dots, n$  be i.i.d. replicates of the vector  $(X, C, Z)'$ , where  $X$  and  $C$  are the survival and right censoring times and  $Z$  is a categorical covariate. It is supposed that  $X_i$  and  $C_i$  are conditionally independent given  $Z_i$ ,  $i = 1, \dots, n$ . We observe the sample  $(T_i, \Delta_i, Z_i)'$ ,  $i = 1, \dots, n$ , where  $T_i = \min\{X_i, C_i\}$  is the observation time and  $\Delta_i = 1_{\{X_i \leq C_i\}}$  is the failure indicator. Let  $F(x|z)$ ,  $x \geq x_0 \geq 0$  and  $F_C(x|z)$ ,  $x \geq x_0$  be the conditional distributions of  $X$  and  $C$ , given  $Z = z$ , respectively. In this paper we address the problem of estimation of the survival function  $S_F(x|z) = 1 - F(x|z)$  when  $x \geq x_0$  is large. The function  $S_F$  is traditionally estimated using the Kaplan-Meier nonparametric estimator (Kaplan and Meier (1958)). Its properties have been extensively studied by numerous authors, including Fleming and Harrington (1991), Andersen, Borgan, Gill and Keiding (1993), Kalbfleisch and Prentice (2002), Klein and Moeschberger (2003). However, in various practical applications, when the time  $x$  is close to or exceeds the largest observed data, the predictions based on the Kaplan-Meier and related estimators are rather uninformative.

For illustration purposes we consider the well known PBC (primary biliary cirrhosis) data from a clinical trial analyzed in Fleming and Harrington (1991). In this trial one observes the censored survival times of two groups of patients: the first one ( $Z = 1$ ) was given the DPCA (D-penicillamine drug) treatment and the second one is the control group ( $Z = 0$ ). The overall censoring rate is about 60%. Here we consider only the group covariate and we are interested to compare the extreme survival probabilities of the patients under study in the two groups. In Figure 1 (left picture) we display the Kaplan-Meier nonparametric curves of the treatment and the control (placebo) groups. From these curves it seems difficult to infer whether the DPCA treatment has an effect on the survival probability. For instance at time  $x = 4745$  (13 years) using the Kaplan-Meier nonparametric estimator (KM), one gets an estimated survival probability  $\hat{S}_{KM}(x|z = 0) = 0.3604$  for

\*First author's affiliation, Université de Bretagne Sud, LMBA, UMR CNRS 6205

<sup>†</sup>Second author's affiliation, Université de Bretagne Sud, LMBA, UMR CNRS 6205

<sup>‡</sup>Third author's affiliation, Université de Bretagne Sud, LMBA, UMR CNRS 6205

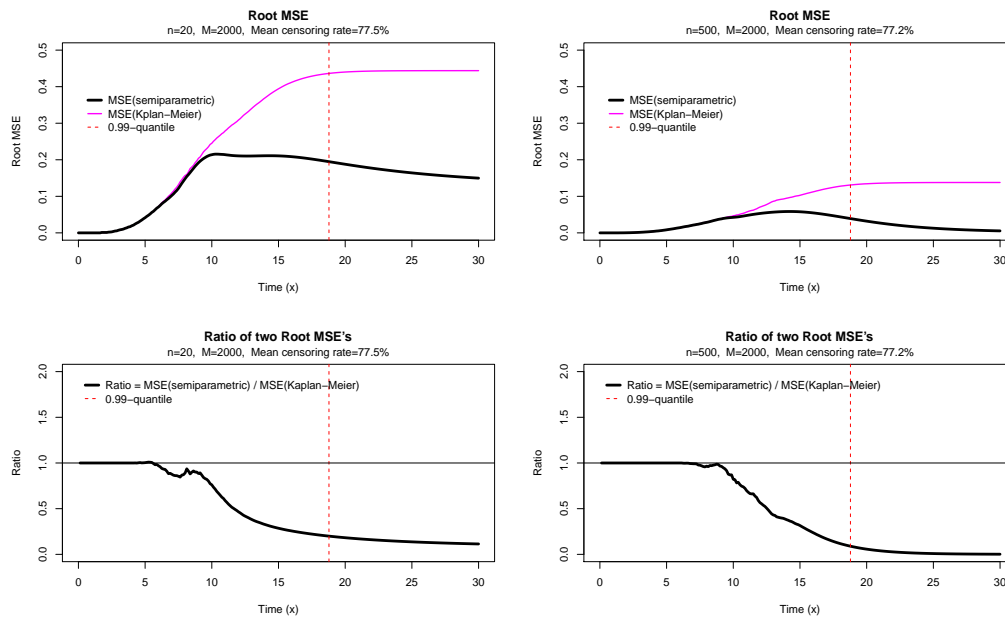


**Figure 1:** We compare two types of prediction of the survival probabilities in DPCA and placebo groups: on the left picture the prediction is based on the Kaplan-Meier estimation and on the right picture the prediction uses a semiparametric approach. The points on the curves correspond to the largest observation time in each group.

the control group and  $\widehat{S}_{KM}(x|z=1) = 0.3186$  for the DPCA treatment group. In this example and in many other applications one has to face the following two drawbacks. First, the estimated survival probabilities  $\widehat{S}_{KM}(x|z)$  are constant for  $x$  beyond the largest (non censored) survival time, which is not quite helpful for prediction purposes. Second, for this particular data set, the Kaplan-Meier estimation suggests that the DPCA treatment group has an estimated long term survival probability slightly lower than that of the control group, which can be explained by the high variability of  $\widehat{S}_{KM}(x|z)$  for large  $x$ . These two points clearly rise the problem of correcting the behavior of the tail of the Kaplan-Meier estimator.

A largely accepted way to estimate the survival probabilities  $S_F(x|z)$  for large  $x$ , is the parametric-based model fitting the hole data starting from the origin. Its advantages are pointed out in Miller (1983), however, it is well known that the bias model can be high if it is misspecified. The more flexible nonparametric Kaplan-Meier estimator would generally be preferred for estimating certain functionals of the survival curve, as argued in Meier, Karrison, Chappell and Xie (2004). In this paper we propose to combine the nonparametric Kaplan-Meier estimator and the parametric-based model into one construction which we call semiparametric Kaplan-Meier estimator (SKM). Our new estimator incorporates a threshold  $t$  in such a way that  $S_F(x|z)$  is estimated by the Kaplan-Meier estimator for  $x \leq t$  and by a parametric-based estimate for  $x > t$ . The main theoretical contribution of the paper is to show that with an appropriate choice of the threshold  $t$  such an estimate is consistent if the tail is correctly specified. In the case when the tail is misspecified we show by simulations that the method is robust. Denote by  $\widehat{S}_t$  the resulting estimator of  $S_F$ , where the parametric-based model is the exponential distribution with mean  $\theta$ . By simulations we have found that  $\widehat{S}_{\widehat{t}}$ , endowed with a data driven threshold  $\widehat{t}$ , outperforms the Kaplan-Meier estimator. As it is seen from Figure 1 (right picture), we obtain at  $x = 4745$  the estimated survival probability  $\widehat{S}_{\widehat{t}_0}(x|z=0) = 0.2739$  for the control group and  $\widehat{S}_{\widehat{t}_1}(x|z=1) = 0.3150$  for the DPCA treatment group, where  $\widehat{t}_0$  and  $\widehat{t}_1$  are the corresponding data driven thresholds. Our predictions are recorded in Table 2 and seem to be more adequate than those based on the Kaplan-Meier estimation. We refer to Section 6, where this example is described in more details.

In Figure 2 we display the root of the mean squared error of the predictions of  $S_F(x|z)$  based on the Kaplan-Meier and the proposed semiparametric Kaplan-Meier estimators as functions of the observation time  $x$ . This is an example where the exponential model for survival and censoring tails is misspecified. The errors are computed within a Monte-Carlo simulation study of size  $M = 2000$  with a gamma distribution modeling the survival and censoring times which do not exhibit exponential behavior in the tail (see Section 5 and Example 2 of Section 2 for details). The advantage of the proposed semiparametric estimator



**Figure 2:** The top line displays simulated root MSE's of the Kaplan-Meier and semiparametric Kaplan-Meier estimators as functions of the time  $x$ . On the bottom line we show the ratio of the two root MSE's displayed on the top line. The two columns correspond to sample sizes  $n = 20$  and  $n = 500$ . The vertical dashed line is the 0.99-quantile of the true distribution of the survival time.

over the Kaplan-Meier estimator can be clearly seen by comparing the two MSE curves. The MSE of the semiparametric estimator is much smaller than that of the Kaplan-Meier estimator for large observation times  $x > q_{0.99}$  but also for mid range observation time values, for example  $x \in [8, q_{0.99}]$ , where  $q_{0.99}$  is the 0.99-quantile of the distribution  $F$ . The proposed extensions of the nonparametric curves are particularly suited for predicting the survival probabilities in the case when the proportion of the censored times is large. This is the case of the mentioned simulated data where the mean censoring rate is about 77%. Note also that we get an improvement over the Kaplan-Meier estimator even for very low sample sizes like  $n = 20$ .

The proposed estimator  $\hat{S}_t$  is sensible to the choice of the threshold  $t$ . The main difficulty is to choose  $t$  small enough, so that the parametric-based part contains enough observation times to ensure a reliable prediction in the tail. At the same time one should choose  $t$  large enough in order to prevent from a large bias due to an inadequate tail fitting. The very important problem of the automatic choice of the threshold  $\hat{t}$  is treated in Section 4, where a procedure which we call *testing-pursuit-selection* is performed in two stages: First we test sequentially the null hypothesis that the proposed parametric-based model fits the data until we detect a chosen alternative. Secondly we select the best model among the accepted ones by penalized model selection. Therefore our testing-pursuit-selection procedure is actually also a goodness-of-fit test for the proposed parametric-based model. The resulting data driven estimator of the tail depends heavily on the testing procedure.

Related result concerning the approximation of the tail by the Cox model (Cox (1972)) can be found in Grama, Tricot and Petiot (2011). The case of continuous multivariate covariate  $Z$  in the context of a Cox model and the use of fitted tails other than the exponential can be treated by similar methods. The models which take into account the cure effects can be reduced to ours after removing the cure fraction. However, these problems are beyond the scope of the paper.

This proceeding is a shortened version of the complete article to appear in "Buletinul

**2. The model and background definitions**

Assume that the survival and right censoring times arise from variables  $X$  and  $C$  which take their values in  $[x_0, \infty)$ , where  $x_0 \geq 0$ . Consider that  $X$  and  $C$  may depend on the categorical covariate  $Z$  with values in the set  $\mathcal{Z} = \{0, \dots, m\}$ . The related conditional distributions  $F(x|z)$  and  $F_C(x|z)$ ,  $x \geq x_0$ , given  $Z = z$ , are supposed to belong to the set  $\mathcal{F}$  of distributions with strictly positive density on  $[x_0, \infty)$ . Let  $f_F(\cdot|z)$  and  $S_F(\cdot|z) = 1 - F(\cdot|z)$  be the conditional density and survival functions of  $X$ , given  $Z = z$ . The corresponding conditional hazard function is  $h_F(\cdot|z) = f_F(\cdot|z)/S_F(\cdot|z)$ , given  $Z = z$ . Similarly,  $C$  has the conditional density  $f_C(\cdot|z)$ , survival function  $S_C(\cdot|z)$  and hazard function  $h_C(\cdot|z) = f_C(\cdot|z)/S_C(\cdot|z)$ , given  $Z = z$ . We also assume the independence between  $X$  and  $C$ , conditionally with respect to  $Z$ . Let the observation time and the failure indicator be

$$T = \min \{X, C\} \quad \text{and} \quad \Delta = 1_{\{X \leq C\}},$$

where  $1_B$  is the indicator function taking the value 1 on the event  $B$  and 0 otherwise. Let  $P_{F,F_C}(dx, d\delta|z)$ ,  $x \in [x_0, \infty)$ ,  $\delta \in \{0, 1\}$  be the conditional distribution of the vector  $\mathbf{Y} = (T, \Delta)'$ , given  $Z = z$ . The density of  $P_{F,F_C}$  is

$$p_{F,F_C}(x, \delta|z) = f_F(x|z)^\delta S_F(x|z)^{1-\delta} f_C(x|z)^{1-\delta} S_C(x|z)^\delta, \tag{2.1}$$

where  $x \in [x_0, \infty)$ ,  $\delta \in \{0, 1\}$ .

Let  $z_i \in \mathcal{Z}$  be the observed value of the covariate  $Z_i$ , where  $Z_i$ ,  $i = 1, \dots, n$  are i.i.d. copies of  $Z$ , and let  $\mathbf{Y}_i = (T_i, \Delta_i)'$ ,  $i = 1, \dots, n$  be a sample of  $n$  vectors, where each vector  $\mathbf{Y}_i$  has the conditional distribution  $P_{F,F_C}(\cdot|z_i)$ , given  $Z_i = z_i$ , for  $i = 1, \dots, n$ . It is clear that, given  $Z = z \in \mathcal{Z}$ , the vectors  $\mathbf{Y}_i$ ,  $i \in \{j : z_j = z\}$  are i.i.d. .

In this paper the problem is to improve the nonparametric Kaplan-Meier estimators of the  $m + 1$  survival probabilities  $S_F(x|z) = 1 - F(x|z)$ ,  $z \in \mathcal{Z}$ , for large values of  $x$ . To this end, we fit the tail of  $F(\cdot|z)$  by the exponential distribution with mean  $\theta > 0$ . Consider the following conditional semiparametric quasi-model

$$F_{\theta,t}(x|z) = \begin{cases} F(x|z), & x \in [x_0, t], \\ 1 - (1 - F(t|z)) \exp(-\frac{x-t}{\theta}), & x > t, \end{cases} \tag{2.2}$$

where  $t \geq x_0$  is a nuisance parameter and  $F(\cdot|z) \in \mathcal{F}$ ,  $z \in \mathcal{Z}$  are functional parameters. The conditional density, survival and hazard functions of  $F_{\theta,t}$  are denoted by  $f_{F_{\theta,t}}$ ,  $S_{F_{\theta,t}}$  and  $h_{F_{\theta,t}}$ , respectively. Note that  $h_{F_{\theta,t}}(x|z) = 1/\theta$ , for  $x > t$ .

We shall focus on two typical models throughout the following examples.

*Example 1 (asymptotically constant hazards).* Consider asymptotically constant survival and censoring hazard functions. This model can be related to the families of distributions in Hall (1982), Hall and Welsh (1984, 1985), Dress (1998) and Grama and Spokoiny (2008) for the extreme value models. Let  $A > 0$ ,  $\theta_{\max} > \theta_{\min} > 0$  be some constants. Consider that the survival time  $X$  has a hazard function  $h_F(\cdot|z)$  such that for some  $\theta_z \in (\theta_{\min}, \theta_{\max})$  and  $\alpha_z > 0$ ,

$$|\theta_z h_F(\theta_z x|z) - 1| \leq A \exp(-\alpha_z x), \quad x \geq x_0. \tag{2.3}$$

Condition (2.3) means that  $h_F(x|z)$  converges to  $\theta_z^{-1}$  exponentially fast as  $x \rightarrow \infty$ . Substituting  $\alpha_z = \alpha'_z \theta_z$ , (2.3) gives  $|h_F(x|z) - \theta_z^{-1}| \leq A' \exp(-\alpha'_z x)$ , where  $A' = A/\theta_{\min}$ .

Similarly, let  $M > 0$ ,  $\gamma_{\max} > \gamma_{\min} > 0$ ,  $\mu > 1$  be some constants. Assume that the hazard function  $h_C(\cdot|z)$  of the censoring time  $C$  satisfies for some  $\gamma_z \in (\gamma_{\min}, \gamma_{\max})$ ,

$$|\theta_z h_C(\theta_z x|z) - \gamma_z| \leq M(1+x)^{-\mu}, \quad x \geq x_0. \tag{2.4}$$

Condition (2.4) is equivalent to saying that  $h_C(x|z)$  approaches  $\gamma_z/\theta_z$  polynomially fast as  $x \rightarrow \infty$ . Substituting  $\gamma_z = \gamma'_z \theta_z$ , (2.4) gives  $|h_C(x|z) - \gamma'_z| \leq M'x^{-\mu}$ , where  $M' = M\theta_{\max}^\mu/\theta_{\min}$ .

For example, conditions (2.3) and (2.4) are satisfied if  $F$  and  $F_C$  coincide with the re-scaled Cauchy distribution  $K_{\mu,\theta}$  defined below. Let  $\xi$  be a variable with the positive Cauchy distribution  $K(x) = 2\pi^{-1} \arctan(x)$ ,  $x \geq 0$ . We define the re-scaled Cauchy distribution by  $K_{\mu,\theta}(x) = 1 - \frac{1-K(\exp((x-\mu)/\theta))}{1-K(\exp(-\mu/\theta))}$ , where  $\mu$  and  $\theta$  are the location and scale parameters. The distribution  $K_{\mu,\theta}$  can be seen as the excess distribution of the variable  $\theta \log \xi + \mu$  over the threshold 0. We leave to the reader the verification that  $K_{\mu,\theta}$  fulfills (2.3) with  $\theta_z = \theta$ ,  $\alpha_z = 2$  and (2.4) with  $\gamma_z = 1$ . The distribution  $K_{\mu,\theta}$  will be used in Section 5 to simulate survival and censoring times.

*Example 2 (non constant hazards).* Now we consider the case when the hazard functions are not asymptotically constant. For instance, this is the case when the survival and censoring times have both gamma distributions. The numerical results presented in Figure 2 and discussed in Section 5 show that the approach of the paper works when conditions (2.3) and (2.4) are not satisfied, which means that the tail probabilities can be estimated by our approach even if the exponential model is misspecified for the tail.

### 3. Consistency of the estimator with fixed threshold

Define the quasi-log-likelihood by  $\mathcal{L}_t(\theta|z) = \sum_{i=1}^n \log p_{F_{\theta,t}, F_C}(T_i, \Delta_i|z_i) 1_{\{z_i=z\}}$ , where  $F_{\theta,t}$  is defined by (2.2) with parameters  $\theta > 0$ ,  $t \geq x_0$  and  $F(\cdot|z) \in \mathcal{F}$ ,  $z \in \mathcal{Z}$ . Taking into account (2.1) and dropping the terms related to the censoring, the partial quasi-log-likelihood is

$$\begin{aligned} \mathcal{L}_t^{\text{part}}(\theta|z) &= \sum_{T_i \leq t, z_i=z} \Delta_i \log h_{F_{\theta,t}}(T_i|z) - \sum_{T_i > t, z_i=z} \Delta_i \log \theta \\ &- \sum_{T_i \leq t, z_i=z} \int_{x_0}^{T_i} h_{F_{\theta,t}}(v|z) dv - \sum_{T_i > t, z_i=z} \left( \int_{x_0}^t h_{F_{\theta,t}}(v) dv + \theta^{-1}(T_i - t) \right), \end{aligned} \tag{3.1}$$

for fixed  $z \in \mathcal{Z}$  and  $t \geq x_0$ . Maximizing  $\mathcal{L}_t^{\text{part}}(\theta|z)$  in  $\theta$ , obviously yields the estimator

$$\hat{\theta}_{z,t} = \frac{\sum_{T_i > t, z_i=z} (T_i - t)}{\hat{n}_{z,t}}, \tag{3.2}$$

where by convention  $0/0 = \infty$  and  $\hat{n}_{z,t} = \sum_{T_i > t, z_i=z} \Delta_i$  is the number of observed survival times beyond the threshold  $t$ .

The estimator of  $S_F(x)$ , for  $x_0 \leq x \leq t$ , is easily obtained by standard nonparametric maximum likelihood approach due to Kiefer and Wolfowitz (1956) (see also Bickel, Klaassen, Ritov and Wellner (1993), Section 7.5). We use the product Kaplan-Meier (KM) estimator (with ties) defined by  $\hat{S}_{KM}(x|z) = \prod_{T_i \leq x} (1 - d_i(z)/r_i(z))$ ,  $x \geq x_0$ , where  $r_i(z) = \sum_{j=1}^n 1_{\{T_j \geq T_i, z_j=z\}}$  is the number of individuals at risk at  $T_i$  and  $d_i(z) = \sum_{j=1}^n 1_{\{T_j=T_i, \Delta_j=1, z_j=z\}}$  is the number of individuals died at  $T_i$  (see Klein and Moeschberger (2003), Section 4.2 and Kalbfleisch and Prentice (2002)). The *semiparametric fixed-threshold Kaplan-Meier estimator* (SFKM) of the survival function takes the

form

$$\widehat{S}_t(x|z) = \begin{cases} \widehat{S}_{KM}(x|z), & x \in [x_0, t], \\ \widehat{S}_{KM}(t|z) \exp\left(-\frac{x-t}{\widehat{\theta}_{z,t}}\right), & x > t, \end{cases} \quad (3.3)$$

where  $\exp\left(-\frac{x-t}{\widehat{\theta}_{z,t}}\right) = 1$  if  $\widehat{\theta}_{z,t} = \infty$ . Similarly, it is possible to use the Nelson-Aalen nonparametric estimator (Nelson 1969, 1972, Aalen, 1976) instead of the Kaplan-Meier one.

Denote by  $n_z = \sum_{i=1}^n 1(z_i = z)$  the number of individuals with profile  $z \in \mathcal{Z}$ . Assume that there is a constant  $\kappa \in (0, 1]$  such that, for any  $z \in \mathcal{Z}$ ,

$$n_z \geq \kappa n. \quad (3.4)$$

Let  $\mathbb{P}$  be the joint distribution of the sample  $\mathbf{Y}_i, i = 1, \dots, n$  and  $\mathbb{E}$  be the expectation with respect to  $\mathbb{P}$ . In the sequel, the notation  $\alpha_n = O_{\mathbb{P}}(\beta_n)$  means that there is a positive constant  $c$  such that  $\mathbb{P}(\alpha_n > c\beta_n, \beta_n < \infty) \rightarrow 0$  as  $n \rightarrow \infty$ , for any two sequences of positive possibly infinite variables  $\alpha_n$  and  $\beta_n$ .

Consider the Kullback-Leibler divergence  $\mathcal{K}(\theta', \theta) = \int \log(dG_{\theta'} / dG_{\theta}) dG_{\theta'}$  between two exponential distributions with means  $\theta'$  and  $\theta$ . By convention,  $\mathcal{K}(\infty, \theta) = \infty$ . It is easy to see that  $\mathcal{K}(\theta', \theta) = \psi(\theta'/\theta - 1)$ , with  $\psi(x) = x - \log(x + 1)$ ,  $x > -1$  and that there are two constants  $c_1$  and  $c_2$  such that  $(\theta'/\theta - 1)^2 \leq c_1 \mathcal{K}(\theta', \theta) \leq c_2 (\theta'/\theta - 1)^2$ , when  $|\theta'/\theta - 1|$  is small enough.

Our main result is the following theorem.

**Theorem 3.1** *Assume condition (3.4). Assume that  $h_F(\cdot|z)$  satisfies (2.3) and  $h_C(\cdot|z)$  satisfies (2.4). Then,*

$$\mathcal{K}(\widehat{\theta}_{z,t_{z,n}}, \theta_z) = O_{\mathbb{P}}\left(\left(\frac{\log n}{n}\right)^{\frac{2\alpha_z}{1+\gamma_z+2\alpha_z}}\right), \quad (3.5)$$

where

$$t_{z,n} = \frac{\theta_z}{1 + \gamma_z + 2\alpha_z} \log n + o(\log n).$$

We give some hints about the optimality of the rate in (3.5). Assume that the survival time  $X$  is exponential, i.e.  $h_F(x|z) = \theta_z^{-1}$  for all  $x \geq x_0$  and  $z \in \mathcal{Z}$ . This ensures that condition (2.3) is satisfied with any  $\alpha > 0$ . Assume conditions (2.4) and (3.4). If there are two constants  $\theta_{\min}$  and  $\theta_{\max}$  such that  $0 < \theta_{\min} \leq \theta_z \leq \theta_{\max} < \infty$ , (3.5) implies  $|\widehat{\theta}_{z,t_{z,n}} - \theta_z| = O_{\mathbb{P}}\left((n^{-1} \log n)^{\frac{\alpha}{1+\gamma_z+2\alpha}}\right)$ , for any  $\alpha > 0$ . This rate becomes arbitrarily close to the  $n^{-1/2}$  rate as  $\alpha \rightarrow \infty$ , since  $\lim_{\alpha \rightarrow \infty} \alpha / (1 + \gamma_z + 2\alpha) \rightarrow 1/2$ . Thus the estimator  $\widehat{\theta}_{z,t_{z,n}}$  almost recovers the usual parametric rate of convergence as  $\alpha$  becomes large whatever is  $\gamma_z > 0$ .

In the case when there are no censoring ( $\gamma_z = 0$ ), after an exponential rescaling our problem can be reduced to that of the estimation of extreme index. If  $\gamma_z \rightarrow 0$  our rate becomes close to  $n^{-\frac{2\alpha_z}{1+2\alpha_z}}$ , which is known to be optimal in the context of the extreme value estimation, see Dress (1998) and Grama and Spokoiny (2008). So our result nearly recovers the best possible rate of convergence in this setting.

#### 4. Testing and automatic selection of the threshold

In this section a procedure of selecting the adaptive estimator  $\widehat{\theta}_z = \widehat{\theta}_{z,\widehat{t}_{z,n}}$  from the family of fixed threshold estimators  $\widehat{\theta}_{z,t}, t \geq x_0$  is proposed. Here the adaptive threshold  $\widehat{t}_{z,n}$

is obtained by a sequential testing procedure followed by a selection using a penalized maximum likelihood. This motivates our condensed terminology *testing-pursuit-selection* used in the sequel. The testing part is actually a multiple goodness-of-fit testing for the proposed parametric-based models, while the threshold  $\hat{t}_{z,n}$  can be seen as a data driven substitute for the theoretical threshold  $t_{z,n}$  defined in Theorem 3.1. For a discussion on the proposed approach we refer the reader to Section 3 of Grama and Spokoiny (2008). In the sequel, for simplicity of notations, we abbreviate  $\hat{t}_z = \hat{t}_{z,n}$ .

Define a semiparametric change-point distribution by

$$F_{\mu,s,\theta,t}(x|z) = \begin{cases} F(x|z), & x \in [x_0, s], \\ 1 - (1 - F(s|z)) \exp\left(-\frac{x-s}{\mu}\right), & x \in (s, t], \\ 1 - (1 - F(s|z)) \exp\left(-\frac{t-s}{\mu}\right) \exp\left(-\frac{x-t}{\theta}\right), & x > t, \end{cases}$$

for  $\mu, \theta > 0$ ,  $x_0 \leq s < t$  and  $F(\cdot|z) \in \mathcal{F}$ . As in Section 3 we find the maximum quasi-likelihood estimators  $\hat{\theta}_{z,t}$  of  $\theta$  and  $\hat{\mu}_{z,s,t}$  of  $\mu$  for fixed  $z \in \mathcal{Z}$  and  $x_0 \leq s < t$ , which are given by (3.2) and

$$\hat{\mu}_{z,s,t} = \frac{\hat{n}_{z,s} \hat{\theta}_{z,s} - \hat{n}_{z,t} \hat{\theta}_{z,t}}{\hat{n}_{z,s,t}},$$

where  $\hat{n}_{z,s,t} = \sum_{s < T_i \leq t, z_i = z} \Delta_i$  and by convention  $0 \cdot \infty = 0$  and  $0/0 = \infty$ .

Consider a constant  $D > 0$ , which will be the critical value in the testing procedure below. Let  $k_0 \geq 3$  be a starting index and  $k_{step}$  be an increment for  $k$ . Let  $\delta', \delta''$  be two positive constants such that  $0 < \delta', \delta'' < 0.5$ . The values  $k_0, k_{step}, \delta', \delta''$  and  $D$  are the parameters of the procedure to be calibrated empirically. Without loss of generality, we consider that the  $T_i$ 's are arranged in the decreasing order:  $T_1 \geq \dots \geq T_n$ . The threshold  $t$  will be chosen in the set  $\{T_1, \dots, T_n\}$ .

The *testing-pursuit-selection* procedure which we propose is performed in two stages. First we test the null hypothesis  $\mathcal{H}_{T_k}(z) : F = F_{\theta, T_k}(\cdot|z)$  against the alternative  $\tilde{\mathcal{H}}_{T_k}(z) : F = F_{\mu, T_k, \theta, T_l}(\cdot|z)$  for some  $\delta'k \leq l \leq (1 - \delta'')k$ , sequentially in  $k = k_0 + ik_{step}$ ,  $i = 0, \dots, [n/k_{step}]$ , until  $\mathcal{H}_{T_k}(z)$  is rejected. Denote by  $\hat{k}_z$  the obtained break index and define the break time  $\hat{s}_z = T_{\hat{k}_z}$ . Second, using  $\hat{k}_z$  and  $\hat{s}_z$  define the adaptive threshold by  $\hat{t}_z = T_{\hat{l}_z}$  with the adaptive index

$$\hat{l}_z = \underset{\delta' \hat{k}_z \leq l \leq (1 - \delta'') \hat{k}_z}{\operatorname{argmax}} \left\{ \mathcal{L}_{T_l}(\hat{\theta}_{z, T_l}|z) - \mathcal{L}_{T_l}(\hat{\theta}_{z, \hat{s}_z}|z) \right\}, \quad (4.1)$$

where the term  $\mathcal{L}_{T_l}(\hat{\theta}_{z, \hat{s}_z}|z)$  is a penalty for getting close to the break time  $\hat{s}_z$ . The resulting adaptive estimator of  $\theta_z$  is defined by  $\hat{\theta}_z = \hat{\theta}_{z, \hat{t}_z}$  and the *semiparametric adaptive-threshold Kaplan-Meier estimator* (SAKM) of the survival function is defined by  $\hat{S}_{\hat{t}_z}(\cdot|z)$ .

For testing  $\mathcal{H}_{T_k}(z)$  against  $\tilde{\mathcal{H}}_{T_k}(z)$  we use the statistic

$$LR_{\max}(T_k|z) = \max_{\delta'k \leq l \leq (1 - \delta'')k} LR(T_k, T_l|z), \quad (4.2)$$

where  $LR(s, t|z)$  is the quasi-likelihood ratio test statistic for testing  $\mathcal{H}_s(z) : F = F_{\theta, s}(\cdot|z)$  against the alternative  $\tilde{\mathcal{H}}_{s,t}(z) : F = F_{\mu, s, \theta, t}(\cdot|z)$ . To compute (4.2), note that by simple calculations, using (3.1) and (3.2),

$$\mathcal{L}_t(\hat{\theta}_{z,t}|z) - \mathcal{L}_t(\theta|z) = \hat{n}_{z,t} \mathcal{K}(\hat{\theta}_{z,t}, \theta), \quad (4.3)$$

where by convention  $0 \cdot \infty = 0$ . Similarly to (4.3), the quasi-likelihood ratio test statistic  $LR(s, t|z)$  is given by

$$LR(s, t|z) = \hat{n}_{z,s,t} \mathcal{K}(\hat{\mu}_{z,s,t}, \hat{\theta}_{z,s}) + \hat{n}_{z,t} \mathcal{K}(\hat{\theta}_{z,t}, \hat{\theta}_{z,s}) \quad (4.4)$$

with the same convention. Note that, by (4.3), the second term in (4.4) can be viewed as the penalized quasi-log-likelihood

$$\begin{aligned} LR_{\text{pen}}(s, t|z) &= \mathcal{L}_t(\hat{\theta}_{z,t}|z) - \mathcal{L}_t(\hat{\theta}_{z,s}|z) \\ &= \hat{n}_{z,t} \mathcal{K}(\hat{\theta}_{z,t}, \hat{\theta}_{z,s}). \end{aligned}$$

Our testing-pursuit-selection procedure reads as follows:

**Step 1.** Set the starting index  $k = k_0$ .

**Step 2.** Compute the test statistic for testing  $\mathcal{H}_{T_k}(z)$  against  $\tilde{\mathcal{H}}_{T_k}(z)$  :

$$LR_{\text{max}}(T_k|z) = \max_{\delta'k \leq l \leq (1-\delta'')k} LR(T_k, T_l|z)$$

**Step 3.** If  $k \leq n - k_{\text{step}}$  and  $LR_{\text{max}}(T_k|z) \leq D$ , increase  $k$  by  $k_{\text{step}}$  and go to Step 2. If  $k > n - k_{\text{step}}$  or  $LR_{\text{max}}(T_k|z) > D$ , let  $\hat{k}_z = k$ ,

$$\hat{l}_z = \underset{\delta'\hat{k}_z \leq l \leq (1-\delta'')\hat{k}_z}{\text{argmax}} LR_{\text{pen}}(T_{\hat{k}_z}, T_l|z),$$

take the adaptive threshold as  $\hat{t}_z = T_{\hat{l}_z}$  and exit.

It may happen that with  $k = k_0$  it holds  $LR_{\text{max}}(T_{k_0}|z) > D$ , which means that the hypothesis that the tail is fitted by the exponential model, starting from  $T_{k_0}$ , is rejected. In this case we resume the procedure with a new augmented  $k_0$ , say with  $k_0$  replaced by  $\lceil \nu_0 k_0 \rceil$ , where  $\nu_0 > 1$ . Finally, if for each such  $k_0$  it holds  $LR_{\text{max}}(T_{k_0}|z) > D$ , we conclude that the tail of the model cannot be fitted with the proposed parametric tail and we estimate the tail by the Kaplan-Meier estimator. Therefore our testing-pursuit-procedure can be seen as well as a goodness-of-fit test for the tail.

Note that the Kullback-Leibler entropy  $\mathcal{K}(\theta', \theta)$  is scale invariant, i.e. satisfies the identity  $\mathcal{K}(\theta', \theta) = \mathcal{K}(\alpha\theta', \alpha\theta)$ , for any  $\alpha > 0$  and  $\theta', \theta > 0$ . Therefore the critical value  $D$  can be determined by Monte Carlo simulations from standard exponential observations. The choice of parameters of the proposed selection procedure is discussed in Section 5.

### 5. Simulation results

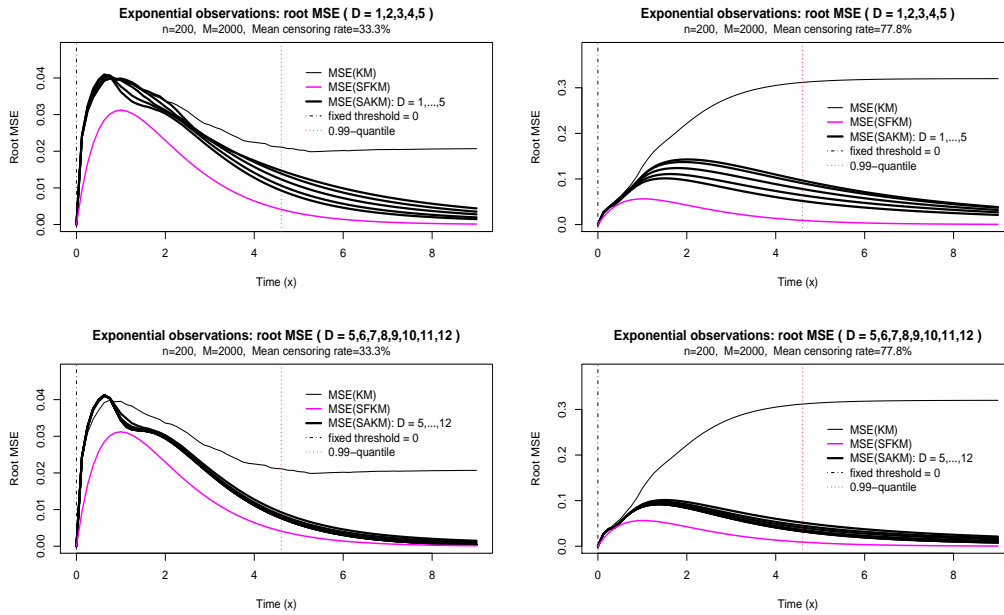
We illustrate the performance of the semiparametric estimator (3.3) with fixed and adaptive thresholds in a simulation study. The survival probabilities  $S_F(x|z)$ , for large values of  $x$ , are of interest.

The mean squared error (MSE) of an estimator  $\hat{S}(\cdot|z)$  of the true survival function  $S_F(\cdot|z)$  is defined by  $MSE_{\hat{S}}(x|z) = \mathbb{E} \left( \hat{S}(x|z) - S_F(x|z) \right)^2$ . The quality of the estimator  $\hat{S}(\cdot|z)$  with respect to the Kaplan-Meier estimator  $\hat{S}_{KM}(\cdot|z)$  is measured by the ratio  $R_{\hat{S}}(x|z) = MSE_{\hat{S}}(x|z) / MSE_{\hat{S}_{KM}}(x|z)$ .

Without loss of generality, we can assume that the covariate  $Z$  takes a fixed value  $z$ . In each study developed below, we perform  $M = 2000$  Monte-Carlo simulations.

We start by giving some hints on the choice of the parameters  $k_0, k_{\text{step}}, \delta', \delta''$  of the testing-pursuit-selection procedure in Section 4. The initial value  $k_0$  controls the variability



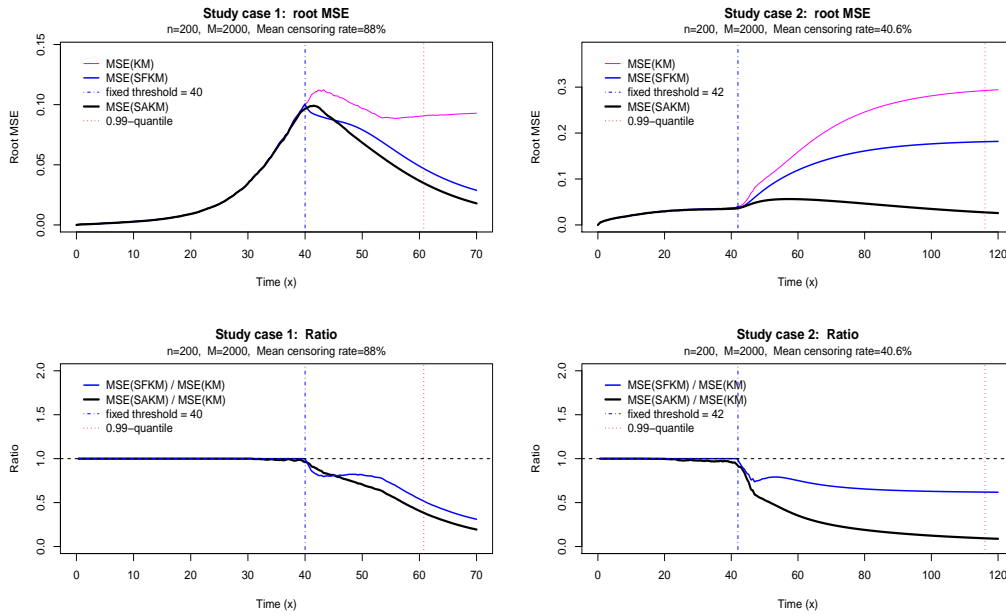


**Figure 3:** The top line displays the root (of type I) MSE's of the Kaplan-Meier estimator (KM), semi-parametric adaptive-threshold Kaplan-Meier estimator (SAKM) with  $D = 1, 2, 3, 4, 5$  and of the exponential model (which coincides with the semiparametric fixed-threshold Kaplan-Meier estimator (SFKM) with threshold fixed at 0). The bottom line displays the same but with  $D = 5, 6, 7, 8, 9, 10, 11, 12$ . The two columns correspond to low and high mean censoring rates.

of the test statistic  $LR_{\max}(T_k|z)$ ,  $k \geq k_0$ . We have fixed  $k_0$  as a proportion of the initial sample size:  $k_0 = n/10$ . The choice  $k_{step} = 5$  is made to speed up the computations. The parameters  $\delta'$  and  $\delta''$  restrict the high variability of the test statistic  $LR(T_k, T_l|z)$  when the change point  $T_l \in [T_k, T_{k_0}]$  is close to the ends of the interval. The values  $\delta' = 0.3$  and  $\delta'' = 0.1$  are retained experimentally. Our simulations show that the adaptive procedure does not depend much on the choice of the parameters  $k_0, k_{step}, \delta', \delta''$ .

To choose the critical value  $D$  we analyze the type I MSE of the SAKM estimator, i.e. the MSE under the null hypothesis that the survival times  $X_1, \dots, X_n$  are i.i.d. standard exponential. We perform two simulations using i.i.d. exponential censoring times  $C_1, \dots, C_n$  with rates 0.5 and 3.5. The size is fixed at  $n = 200$ , but the results are quite similar for other sizes. The root MSE's as functions of the time  $x$  are given in Figure 3. For comparison, in Figure 3 we also included the MSE's corresponding to the parametric-based exponential modeling which coincides with the SFKM estimator having the threshold fixed at 0. Note that the MSE's calculated when the critical values are  $D = 1, 2, 3, 4, 5$ , decrease as  $D$  increases (see the top displays), while for  $D = 5, 6, 7, 8, 9, 10, 11, 12$  the MSE's almost do not depend on  $D$  (see the top and bottom displays). The simulations show that the type I MSE decreases as  $D$  increases and stabilizes for  $D \geq 5$ . From these plots we conclude that the limits for the critical value  $D$  can be set between  $D_0 = 5$  and  $D_1 = 7$  without important loss in the type I MSE.

It is interesting to note that the adaptive threshold  $\hat{t}_z$  is relatively stable to changes of  $D$ . A typical trajectory of the test statistic  $LR_{\max}(T_k|z)$  as function of  $T_k$  is drawn in Figure 6 (left). Despite the fact that the break time  $\hat{s}_z = T_{\hat{k}_z}$  strongly depends on the critical value  $D$  (in the picture  $D = 5.8$ ), we found that the adaptive threshold  $\hat{t}_z = T_{\hat{l}_z}$ , which maximizes the penalized quasi-log-likelihood  $LR_{\text{pen}}(T_{\hat{k}_z}, T_l|z)$  in Figure 6 (right), is stable to the local changes of the break time  $\hat{s}_z = T_{\hat{k}_z}$  and thus is also quasi stable to relatively small changes of  $D$ .



**Figure 4:** The top line displays the root (of type II) MSE's of three estimators:  $\widehat{S}_{KM}$  (KM),  $\widehat{S}_t$  (SFKM) and  $\widehat{S}_{\widehat{t}_z}$  (SAKM). The bottom line displays the corresponding ratios of the root MSE's on the top line. The critical value  $D$  in the SAKM is set to 6.

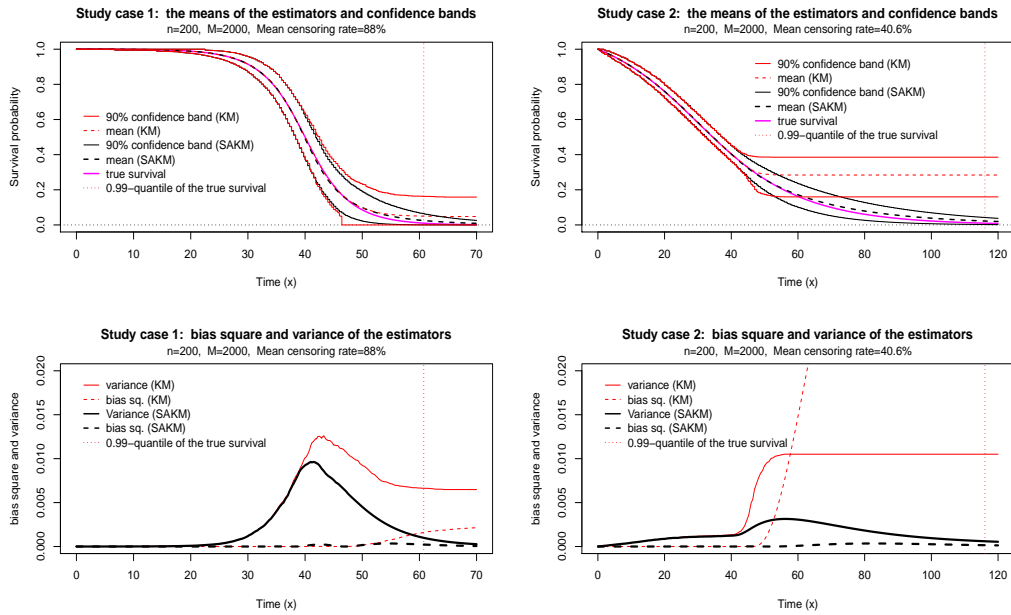
For our simulations we fix the value  $D = 6$ . Below we give some evidence that the SAKM estimator with this critical value has a reasonable type II MSE, under the hypothesis that the  $X_i$ 's have a distribution  $F$  alternative to the standard exponential. Our simulations show that the type II MSE's are quite similar for several families we have tested. We have chosen the following two typical cases which are representative for all these families.

**Study case 1 (low tail censoring rate).** We generate a sequence of  $n = 200$  i.i.d. survival times  $X_i, i = 1, \dots, n$  from the re-scaled Cauchy distribution  $K_{\mu_X, \theta_X}$  with location parameter  $\mu_X = 40$  and scale parameter  $\theta_X = 5$  (see Section 2). The censoring times  $C_i, i = 1, \dots, n$  are i.i.d. from the re-scaled Cauchy distribution  $K_{\mu_C, \theta_C}$  with location parameter  $\mu_C = \mu_X - 20 = 20$  and scale parameter  $\theta_C = 2\theta_X = 10$ . The (overall) mean censoring rate in this example is about 88%. However, the censoring rate for high observation times is about 33%.

**Study case 2 (high tail censoring rate).** As before, let us fix  $n = 200$ . The  $X_i$ 's,  $i = 1, \dots, n$  are i.i.d. from  $K_{\mu_X, \theta_X}$  with  $\mu_X = 30$  and  $\theta_X = 20$ . The  $C_i$ 's,  $i = 1, \dots, n$  are i.i.d. from  $K_{\mu_C, \theta_C}$  with  $\mu_C = \mu_X + 10 = 40$  and  $\theta_C = \theta_X/10 = 2$ . The (overall) mean censoring rate in this example is about 40%, however, the censoring rate among the high observation times is about 91%.

We evaluate the performance of the SFKM and SAKM estimators  $\widehat{S}_t(x|z)$  and  $\widehat{S}_{\widehat{t}_z}(x|z)$  with respect to the KM estimator  $\widehat{S}_{KM}(x|z)$ . In Figure 4 we display the root  $MSE_{\widehat{S}}(x|z)$  and the ratio  $R_{\widehat{S}}(x|z)$  for the three estimators as functions of the time  $x$ . From these plots we can see that both root  $MSE_{\widehat{S}_t}(x|z)$  and root  $MSE_{\widehat{S}_{\widehat{t}_z}}(x|z)$  are equal to the root  $MSE_{\widehat{S}_{KM}}(x|z)$  for small values of  $x$  and become smaller for large values of  $x$ , which shows that the SFKM and SAKM estimators improve the KM estimator.

In Figure 5 (top displays), for each fixed  $x$ , we show the confidence bands containing 90% of the values of  $\widehat{S}_{KM}(x|z)$  and  $\widehat{S}_{\widehat{t}_z}(x|z)$ . From these plots we see the ability of the model to fit the data and at the same time to give satisfactory predictions. Compared to those provided by the KM estimator which predicts a constant survival probability for large  $x$ , our predictions are more realistic.



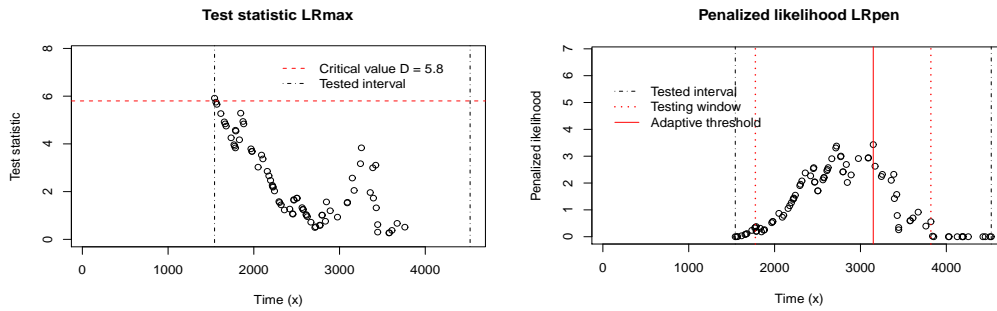
**Figure 5:** The top line displays the true survival  $S_F$  and the estimated means of  $\hat{S}_{KM}(\cdot|z)$  (KM) and  $\hat{S}_{t_z}(\cdot|z)$  (SAKM). We give confidence bands containing 90% of the trajectories for each fixed time  $x$ . The bottom line displays the corresponding biases square and variances.

**Table 1:** Simulations with gamma distributions for survival and censoring times

$x$	5	6	7	8	9	10	11	12	13
$S_F(x z)$	0.9682	0.9161	0.8305	0.7166	0.5874	0.4579	0.3405	0.2424	0.1658
Mean of $\hat{S}_{t_z}(x z)$	0.9679	0.9159	0.8318	0.7107	0.5686	0.4504	0.3575	0.2853	0.2287
Mean of $\hat{S}_{KM}(x z)$	0.9679	0.9159	0.8306	0.7160	0.5875	0.4581	0.3399	0.2472	0.1888
Root $MSE_{\hat{S}_{t_z}}(x z)$	0.0135	0.0225	0.0336	0.0461	0.0552	0.0606	0.0702	0.0831	0.0940
Root $MSE_{\hat{S}_{KM}}(x z)$	0.0135	0.0225	0.0345	0.0466	0.0604	0.0758	0.0933	0.1144	0.1284
$x$	14	15	16	17	18	19	20	21	22
$S_F(x z)$	0.1094	0.0699	0.0433	0.0261	0.0154	0.0089	0.0050	0.0028	0.0015
Mean of $\hat{S}_{t_z}(x z)$	0.1841	0.1487	0.1205	0.0979	0.0798	0.0652	0.0534	0.0439	0.0361
Mean of $\hat{S}_{KM}(x z)$	0.1586	0.1453	0.1411	0.1403	0.1402	0.1402	0.1402	0.1402	0.1402
Root $MSE_{\hat{S}_{t_z}}(x z)$	0.0997	0.0998	0.0952	0.0876	0.0785	0.0690	0.0599	0.0515	0.0441
Root $MSE_{\hat{S}_{KM}}(x z)$	0.1384	0.1503	0.1627	0.1731	0.1804	0.1850	0.1877	0.1893	0.1902

In Figure 5 (bottom displays) we show the bias square and the variance of  $\hat{S}_{KM}(\cdot|z)$  and  $\hat{S}_{t_z}(\cdot|z)$ . From these plots we see that the variance of  $\hat{S}_{t_z}(\cdot|z)$  is smaller than that of  $\hat{S}_{KM}(\cdot|z)$  in the two study cases. We conclude the same for their biases. However, the bias of  $\hat{S}_{KM}(\cdot|z)$  is large in the study case 2 (right bottom display) because of a high censoring rate in the tail.

**The case of non constant hazards (see Example 2 of Section 2).** The previous study is performed for models satisfying conditions (2.3) and (2.4). Now we consider the case when these conditions are not satisfied. Let  $X$  and  $C$  be generated from gamma distributions whose hazard rate function can be easily verified not to be asymptotically constant (in fact it is slowly varying at infinity). The survival time  $X$  is gamma with shape parameter 10 and rate parameter 1 and the censoring time  $C$  is gamma with shape parameter 8.5 and rate parameter 1.2. The mean censoring rate in this example is about 77%. The results of the simulations are given in Figure 2 ( $n = 20$  left picture,  $n = 500$  right picture) and Table 1 ( $n = 500$ ) for  $\hat{S}_{KM}(\cdot|z)$  and  $\hat{S}_{t_z}(\cdot|z)$ . They show that for these distributions the SAKM estimator gives a smaller root MSE than the KM estimator even when the sample size is low ( $n = 20$ ) and  $x$  is in the range of the data.



**Figure 6:** For the placebo group of PBC data we display the test statistics  $LR_{\max}(T_k|z)$  as function of  $T_k$  (left) and  $LR_{\text{pen}}(T_{\hat{k}}, T_l|z)$  as function of  $T_l$  (right). The tested interval and the testing window are given by  $[T_{\hat{k}}, T_{k_0}]$  and  $[T_{(1-\delta'')\hat{k}}, T_{\delta'\hat{k}}]$  respectively. The critical value  $D$  is fixed to 5.8.

**Table 2:** Predicted survival probabilities for PBC data

$x$ : years	3	4	5	6	7	8	9	10	11
$x$ : days	1095	1460	1825	2190	2555	2920	3285	3650	4015
DPCA: KM	0.8256	0.7635	0.7077	0.6613	0.5842	0.5417	0.4778	0.4247	0.4247
DPCA: SAKM	0.8256	0.7635	0.7077	0.6595	0.5934	0.5340	0.4805	0.4323	0.3890
Placebo: KM	0.7911	0.7398	0.7146	0.6950	0.6566	0.6055	0.5461	0.4563	0.3604
Placebo: SAKM	0.7911	0.7398	0.7146	0.6950	0.6566	0.6055	0.5497	0.4619	0.3881
$x$ : years	12	13	14	15	16	17	18	19	20
$x$ : days	4380	4745	5110	5475	5840	6205	6570	6935	7300
DPCA: KM	0.3186	0.3186	0.3186	0.3186	0.3186	0.3186	0.3186	0.3186	0.3186
DPCA: SAKM	0.3501	0.3150	0.2834	0.2550	0.2295	0.2065	0.1858	0.1672	0.1505
Placebo: KM	0.3604	0.3604	0.3604	0.3604	0.3604	0.3604	0.3604	0.3604	0.3604
Placebo: SAKM	0.3260	0.2739	0.2302	0.1934	0.1625	0.1365	0.1147	0.0964	0.0810

## 6. Application to real data

As an illustration we deal with the well known randomized trial in primary biliary cirrhosis (PBC) from Fleming and Harrington (1991) (see Appendix D.1). PBC is a rare but fatal chronic liver disease and the analyzed event is the patient’s death. The trial was open for patient registration between January 1974 and May 1984. The observations lasted until July 1986, when the disease and survival status of the patients were recorded. There were  $n = 312$  patients registered for the clinical trial, including 125 patients who died. The censored times were recorded either for patients which had been lost to follow up or had undergone liver transplantation or was still alive at the study analysis time (July 1986). The number of censored times is 187 and the censoring rate is about 59.9%. The last observed time is 4556 which is a censored time. Ties occur for the following three times: 264, 1191 and 1690. So there are 122 separate times for which we can observe at least one event. Two treatment groups of patients were compared: the first one ( $Z = 1$ ) of size  $n_1 = 158$  was given the DPCA (D-penicillamine drug). The second group ( $Z = 0$ ) of size  $n_0 = 154$  was the control (placebo) group. In this example we consider only the group covariate. We are interested to predict the survival probabilities of the patients under study in both groups.

The survival curves based on the KM and SAKM estimators for each group are displayed in Figure 1. The numerical results on the predictions appear in Table 2. In this table, the time is running from 3 years ( $x = 1095$  days) up to 20 years ( $x = 7300$ ) with the step 1 year equivalent to 365 days for convenience.

Based on the usual KM estimator, the following two conclusions can be made: A1) The constant predictions for extreme survival probabilities in both groups appear to be too optimistic after the largest (non censored) survival time. B1) The DPCA treatment appears to be less efficient than placebo in the long term. The statistical analysis with the SAKM estimator leads to more realistic conclusions: A2) The survival probabilities of each group

extrapolate the tendency of the KM estimator as the time is increasing, and B2) the DPCA treatment is more efficient than placebo. For example, from the results in the Table 2 we obtain that the survival probability in 20 years is about 2 times higher for the DPCA group than for the control group.

From the left picture of Figure 6 we see that the test statistic  $LR_{\max}(T_k|0)$  for the control group ( $Z = 0$ ) reaches the critical value  $D = 5.8 \in [D_0, D_1]$  for  $k = \hat{k}_0 = 90$ . Thus the hypotheses  $\mathcal{H}_{s_0}(0)$  was rejected for the break time  $\hat{s}_0 = T_{\hat{k}_0} = 1542$ . The adaptive threshold  $\hat{t}_0$  is chosen via the maximization of the penalized quasi-log-likelihood (4.1). In the right picture of Figure 6 we see that the maximum is attained for the adaptive index  $\hat{l}_0 = 30$  and threshold  $\hat{t}_0 = T_{\hat{l}_0} = 3149$ . Thus, our testing-pursuit-selection procedure has captured the "convex bump" on the control Kaplan-Meier curve (for  $Z = 0$ ) between the times 2000 and 3500, which is easily seen in the right picture of Figure 1.

## 7. Conclusion

This article deals with estimation of the survival probability in the framework of censored survival data. While the Kaplan-Meier estimator provides a flexible estimate of the survival function in the range of the data it can be improved for prediction of the extreme values, especially when the censoring rate in the tail is high. We propose a new approach based on the Kaplan-Meier estimator by adjusting a parametric correction to the tail beyond a given threshold  $t$ .

In applications the threshold  $t$  usually is not known. To overcome this we propose a testing-pursuit-selection procedure which yields an adaptive threshold  $t = \hat{t}_{z,n}$  in two stages: a sequential hypothesis testing and an adaptive choice of the threshold based on the maximization of a penalized quasi-log-likelihood. This testing-pursuit-selection procedure provides also a goodness-of-fit test for the parametric-based part of the model.

We perform numerical simulations with both the fixed and adaptive threshold estimators. Our simulations show that both estimators improve the Kaplan-Meier estimator not only in the long term, but also in a mid range inside the data. Comparing the fixed threshold and adaptive threshold estimators, we found that the adaptive choice of the threshold significantly improves on the quality of the predictions of the survival function. The improvement over the Kaplan-Meier estimator is especially effective when the censoring rate in the tail is high.

## REFERENCES

- Aalen, O. O. (1976), "Nonparametric Inference in Connection with Multiple Decrements Models," *Scand. J. Statist.*, 3, 15–27.
- Andersen, P. K., Borgan, Ø., Gill, R. D. and Keiding, N. (1993), *Statistical Models Based on Counting Processes*, New York: Springer-Verlag.
- Bickel, P. J., Klaassen, C. A., Ritov, Y. and Wellner, J. A. (1993), *Efficient and Adaptive Estimation for Semiparametric Models*, The Johns Hopkins University Press.
- Cox, D. R. (1972), "Regression Models and Life Tables," *J. Roy. Statist. Soc.*, Ser. B, 34, 187–220.
- Dress, H. (1998), "Optimal Rates of Convergence for Estimates of the Extreme Value Index," *Ann. Statist.*, 26, 434–448.
- Fleming, T. and Harrington, D. (1991), *Counting Processes and Survival Analysis*, Wiley.
- Grama, I. and Spokoiny, V. (2008), "Statistics of Extremes by Oracle Estimation," *Ann. Statist.*, 36, 1619–1648.
- Grama, I., Tricot, J.-M. and Petiot, J.-F. (2011), "Estimation of Survival Probabilities by Adjusting a Cox Model to the Tail," *C.R. Acad. Sci. Paris, Ser. I.*, 349, 807–811.
- Hall, P. (1982), "On Some Simple Estimates of an Exponent of Regular Variation," *J. Roy. Statist. Soc.*, Ser. B, 44, 37–42.
- Hall, P. and Welsh, A. H. (1984), "Best Attainable Rates of Convergence for Estimates of Parameters of Regular Variation," *Ann. Statist.*, 12, 1079–1084.

- Hall, P. and Welsh, A. H. (1985), "Adaptive Estimates of Regular Variation," *Ann. Statist.*, 13, 331–341.
- Kalbfleisch, J. D. and Prentice, R. L. (2002), *The Statistical Analysis of Failure Time Data*, Wiley.
- Kaplan, E. I. and Meier, P. (1958), "Nonparametric Estimation from Incomplete Observation," *J. Amer. Statist. Assoc.*, 53, 457–481.
- Kiefer, J. and Wolfowitz, J. (1956), "Consistency of the Maximum Likelihood Estimator in the Presence of Infinitely Many Nuisance Parameters," *Ann. Math. Statist.*, 27, 887–906.
- Klein, J. P. and Moeschberger, M. L. (2003), *Survival Analysis: Techniques for Censored and Truncated Data*, Springer.
- Meier, P., Karrison, T., Chappell, R. and Xie, H. (2004), "The Price of Kaplan-Meier," *J. Amer. Statist. Assoc.*, 99, 890–896.
- Miller, R. (1983), "What Price Kaplan-Meier ? ." *Biometrics*, 39, 1077-1081.
- Nelson, W. B. (1969), "Hazard Plotting for Incomplete Failure Data," *J. Qual. Technol.*, 1, 27–52.
- Nelson, W. B. (1972), "Theory and Applications of Hazard Plotting for Censored Failure Data," *Technometrics*, 14, 945–965.