

# Adaptive estimation of hazard functions

Sebastian Döhler and Ludger Rüschendorf

*University of Freiburg*

## Abstract

In this paper we obtain convergence rates for sieved maximum-likelihood estimators of the log-hazard function in a censoring model. We also establish convergence results for an adaptive version of the estimator based on the method of structural risk-minimization. Applications are discussed to tensor product spline estimators as well as to neural net and radial basis function sieves. We obtain simplified bounds in comparison to the known literature. This allows to derive several new classes of estimators and to obtain improved estimation rates. Our results extend to a more general class of estimation problems and estimation methods (minimum contrast estimators).

*Keywords:* adaptive estimation, sieved maximum likelihood, neural nets, structural risk minimization, hazard functions

## 1 Introduction

In this paper we establish convergence rates for sieved maximum-likelihood estimators for the log-hazard function in a censoring model. We also establish an adaptive version of the estimator based on the method of structural risk minimization (complexity regularization) as introduced in Vapnik (1995). Our results are obtained for general sieves and then are applied to some special types of sieves like tensor product splines or neural nets. We also state extensions of these results to more general estimation procedures (minimum contrast estimators) and to other types of estimation problems like regression problems comparable to those considered in Birgé and Massart (1998) or in Barron, Birgé, and Massart (1999). For related results see also Krzyzak and Linder (1998), Lugosi and Zeger (1995), Wong and Shen (1995), Yang and Barron (1998), and Kohler(1999a,b).

Sieved ML-estimators are defined in the general framework of empirical risk minimization. The main tools for their analysis are from empirical process theory. The main part of the proof of convergence properties is to establish an exponential maximal inequality for the log-likelihood functional

and to obtain estimates for the covering numbers and Vapnik-Cervonenkis dimension of the involved function classes. In comparison to a similar maximal inequality in Birgé and Massart (1998) we avoid the somewhat complicated condition M2 on control of fluctuations in the  $L_\infty$ -metric and replace it by some more handy growth condition on  $L^1$ -covering numbers. Our  $L^1$ -covering condition is related to condition  $M_{1,[\cdot]}$  ( $L^1$ -metric with bracketing) in Barron, Birgé, and Massart (1999) which is used in that paper to deal with model selection in a general framework and applied to several examples (see sections 4.1.5 and 4.1.6). In comparison our covering condition seems to be particularly simple and well suited for the examples considered in this paper. Our proof is based on an exponential maximal inequality in Lee, Bartlett, and Williamson (1996). In several examples we obtain improved convergence rates in comparison to the literature and some of them are established for the first time in this paper.

In the case of tensor product splines we obtain up to a logarithmic factor the optimal convergence rate in the minimax sense in smoothness classes as derived in Kooperberg, Stone, and Truong (1995b) the only paper on convergence rates in this context so far. For general background on censoring models and reference to martingale based estimation methods we refer to Andersen, Borgan, Gill, and Keiding (1993). Related consistency results for kernel type estimators and further references on nonparametric functional estimation of hazard functions can be found in van Keilegom and Veraverbeke (2001). In comparison to Kooperberg, Stone, and Truong (1995b) we consider the stronger MISE (mean integrated square error). The convergence rate obtained in this paper depends on the smoothness parameter  $p$  of the underlying class of hazard functions as well as on the dimension of the covariables. Some empirical study of an adaptive estimator ('HARE') has been given in Kooperberg, Stone, and Truong (1995a). The related complexity regularized estimator introduced in section 4 of this paper is proved to be adaptive up to a logarithmic order and, therefore, approximatively minimax adaptive. We also discuss applications to general net sieves assuming that the log hazard function allows an integral representation. In particular we consider neural nets, radial basis-function nets and wavelet nets. For further details related to this paper we refer to the dissertation of Döhler (2000b). Some related consistency results (without rates) have been given in Döhler (2000a).

The paper is organized as follows: In chapter 2 we establish an exponential inequality for the log-likelihood functional in the case of right censored data and indicate how similar exponential inequalities can be derived in a general framework. We use this result to obtain general error bounds for sieved ML-estimators (chapter 3) and their complexity regularized versions (chapter 4). In chapter 5 we apply these results to tensorproduct splines and neural net type sieves. We conclude the paper with a short outlook.

The framework of hazard function estimation is as in Kooperberg, Stone,

and Truong (1995b) where however also additive models are considered. Let  $(\Omega, \mathcal{A}, P)$  be the underlying probability space,  $T : \Omega \rightarrow \mathbb{R}_+$  a survival (failure) time,  $C : \Omega \rightarrow \mathcal{T}$  a bounded censoring time,  $X : \Omega \rightarrow \mathcal{X} = [0, 1]^k$  a vector of covariates, and  $Y = T \wedge C$  the observable time. By normalization we assume without loss of generality that  $\mathcal{T} = [0, 1]$ . With the censoring indicator  $\delta = 1_{(T \leq C)}$  (right censoring) the observation vector is  $Z = (X, Y, \delta)$ . We assume existence of a conditional density  $f_0(t|x)$  and denote by  $F_0(t|x)$  the conditional distribution function of  $T$  given  $X = x$ . Further we define  $\lambda_0(t|x) = \frac{f_0(t|x)}{\bar{F}_0(t|x)}$  the conditional hazard function, with conditional survival function  $\bar{F}_0(t|x) = 1 - F_0(t|x)$ , and finally  $\alpha_0(t|x) = \log \lambda_0(t|x)$  the conditional log-hazard function. Based on iid data  $(T_1, C_1, X_1), \dots, (T_n, C_n, X_n)$  respectively the corresponding observed data  $Z_i = (X_i, Y_i, \delta_i)$ ,  $1 \leq i \leq n$ , our aim is to estimate the underlying conditional log-hazard function  $\alpha_0$ .

According to Kooperberg, Stone, and Truong (1995b) the conditional log-likelihood of a sample  $z_1, \dots, z_n$  is given by

$$L_n(\alpha) = \sum_{i=1}^n \ell(z_i, \alpha) \quad (1.1)$$

where  $\ell((x, y, \delta), \alpha) = \delta \alpha(y, x) - \int_0^y \exp \alpha(u, x) du$ .

The underlying log-hazard function is assumed to be in a class  $\mathcal{F}$  of functions on  $\mathcal{T} \times \mathcal{X}$  to be specified later. Generally we assume that  $\alpha$  is bounded on  $\mathcal{T} \times \mathcal{X}$  and that  $T$  and  $C$  are conditionally independent given  $X$ .

Let

$$\Lambda(\alpha) = EL_1(\alpha) \quad (1.2)$$

denote the expected conditional log-likelihood function. Then  $\Lambda$  is maximized at the underlying conditional log-hazard functional  $\alpha_0$ . The sieved maximum-likelihood estimator  $\hat{\alpha}_n$  will be defined by

$$\hat{\alpha}_n = \arg \max_{\alpha \in \mathcal{F}_n} L_n(\alpha) \quad (1.3)$$

over some net (sieve)  $\mathcal{F}_n \subset \mathcal{F}$  depending on the number  $n$  of observations.

For the ‘ $\Lambda$ -distance’ between an arbitrary element  $\alpha \in \mathcal{F}$  and the underlying true  $\alpha_0$  the following representation is useful (see Döhler (2000a)):

$$\begin{aligned} |\Lambda(\alpha) - \Lambda(\alpha_0)| &= \Lambda(\alpha_0) - \Lambda(\alpha) \\ &= \int_{\mathcal{T} \times \mathcal{X}} \bar{F}_{C|X} G(\alpha - \alpha_0) dP^{(T, X)} \end{aligned} \quad (1.4)$$

where  $\bar{F}_{C|X}$  is the conditional survival-function of the censoring time  $C$  and  $G(y) = \exp(y) - (1 + y)$ . A standard argument leads to the decomposition of the estimation error of the ML-estimator  $\hat{\alpha}_n$  (in  $\Lambda$ -distance) in an

approximation error and a stochastic error:

$$|\Lambda(\hat{\alpha}_n) - \Lambda(\alpha_0)| \leq \inf_{\alpha \in \mathcal{F}_n} |\Lambda(\alpha) - \Lambda(\alpha_0)| + 2 \sup_{\alpha \in \mathcal{F}_n} \left| \frac{1}{n} L_n(\alpha) - \Lambda(\alpha) \right|. \quad (1.5)$$

The main tool for proving convergence rates for the stochastic error of  $\hat{\alpha}_n$  will be an exponential maximal inequality derived in section 2. As in Kooperberg, Stone, and Truong (1995b) we introduce the  $L^p$ -distance on  $\mathcal{F}$  modified by the conditional survival function:

$$\|\alpha - \beta\|_p^p = \int_{T \times \mathcal{X}} \bar{F}_{C|X} |\alpha - \beta|^p dP^{(T,X)}. \quad (1.6)$$

From the representation in (1.4) one obtains (see Döhler (2000a)):

$$\Lambda\text{-convergence of } \alpha_n \rightarrow \alpha_0 \text{ implies } \|\alpha_n - \alpha_0\|_1 \rightarrow 0. \quad (1.7)$$

Also for  $\alpha, \beta \in \mathcal{F}$ ,  $|\alpha| \leq M$ ,  $|\beta| \leq M$  holds:

$$k \|\alpha - \alpha_0\|_2^2 \leq |\Lambda(\alpha) - \Lambda(\beta)| \leq k' \|\alpha - \beta\|_2^2 \quad (1.8)$$

where  $k = k(M) = \frac{1}{4M}$ ,  $k' = k'(M) = \frac{\exp(2M)}{4M^2}$ .

For the proof of (1.8) define

$$F(y) = \begin{cases} \frac{G(y)}{y^2} & \text{if } y \neq 0, \\ \frac{1}{2} & \text{if } y = 0, \end{cases}$$

where  $G$  is as in (1.4). Then it is easy to establish that  $F$  is strictly increasing on  $\mathbb{R}$  and  $F(2M) \leq k'(M)$ ,  $F(-2M) \geq k(M)$  for  $M \geq 1$ . Therefore,  $k(M)y^2 \leq G(y) \leq k'(M)y^2$  which implies  $k(M)(\alpha - \beta)^2 \leq G(\alpha - \beta) \leq k'(M)(\alpha - \beta)^2$  and the result follows.

Finally we note that for  $\beta, \alpha \in \mathcal{F}$ ,  $|\beta|, |\alpha| \leq M$

$$E(\ell(Z, \alpha) - \ell(Z, \beta))^2 \leq (B_0 \exp M + 1)^2 \|\alpha - \beta\|_2^2 \quad (1.9)$$

where  $B_0 = \exp M \exp(\exp M)$ . For the proof see Döhler (2000b, Proposition 2.9). So the  $L_2$ -norm allows to control the expected squared loss in the likelihood.

## 2 Exponential maximal inequality for the log-likelihood functional

In this section we derive an exponential maximal inequality for the log-likelihood functional  $L_n(\alpha)$ . The proof is based on the following exponential inequality of Lee, Bartlett, and Williamson (1996) which was used in their

paper and also in Kohler (1997) and Krzyzak and Linder (1998) for regression estimation by minimum  $L^2$ -empirical risk estimators. Based on the error decomposition in (1.5) and relations (1.6), (1.7), and (1.8) we will apply this result to obtain convergence rates of ML-estimators for right censored data.

Let  $N(\varepsilon, \mathcal{F}, d)$  denote the  $\varepsilon$ -covering number of  $\mathcal{F}$  with respect to a metric  $d$ . In the following we will use  $L^p$ -metrics denoted by  $d_{L^p(\mu)}$  on certain  $L^p$ -spaces. The notion of permissibility of  $\mathcal{F}$  denotes a weak measurability condition on  $\mathcal{F}$  allowing to measure sets involving suprema over  $f \in \mathcal{F}$ . For a formal definition see (Pollard 1984, pg. 196). For this and related notions and some basic results on VC-classes we refer to van der Vaart and Wellner (1996) and Pollard (1990).

**Theorem 2.1 (Lee, Bartlett, and Williamson (1996))**

Let  $\mathcal{F}$  be a permissible class of functions on  $\mathcal{Z}$  with  $|f| \leq K_1$ ,  $Ef \geq 0$  and  $Ef^2 \leq K_2Ef$  for all  $f \in \mathcal{F}$ . Let  $\nu, \nu_c > 0$ ,  $0 < \alpha \leq \frac{1}{2}$ , then for  $m \geq \max\{4(K_1 + K_2)/\alpha^2(\nu + \nu_c), K_1^2/\alpha^2(\nu + \nu_c)\}$  holds:

$$\begin{aligned} P \left( \sup_{f \in \mathcal{F}} \frac{Ef - \frac{1}{m} \sum_{i=1}^m f(z_i)}{\nu + \nu_c + Ef} \geq \alpha \right) \\ \leq \sup_{\bar{z} \in \mathcal{Z}^{2m}} 2N \left( \frac{\alpha\nu_c}{4}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})} \right) \exp \left( \frac{-3\alpha^2\nu m}{4K_1 + 162K_2} \right) \\ + \sup_{\bar{z} \in \mathcal{Z}^{2m}} 4N \left( \frac{\alpha\nu_c}{4K_1}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})} \right) \exp \left( \frac{-\alpha^2\nu m}{2K_1^2} \right) \end{aligned} \quad (2.1)$$

where  $\nu_{\bar{z}} = \frac{1}{2m} \sum_{i=1}^{2m} \delta_{\bar{z}_i}$ .

Let now  $\mathcal{Z} = \mathcal{X} \times \mathcal{T} \times \{0, 1\}$  and for  $z_i = (x_i, y_i, \delta_i) \in \mathcal{Z}$ ,  $1 \leq i \leq n$  and  $\bar{z} = (z_1, \dots, z_n)$  let  $\nu_{\bar{z}} = \frac{1}{n} \sum_{i=1}^n \delta_{z_i}$ ,  $\tilde{\nu}_{\bar{z}} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  and let  $U[0, 1]$  be the uniform distribution on  $[0, 1]$ .

**Theorem 2.2 (maximal inequality for the log-likelihood)**

There exists  $B_0 = B_0(\|\alpha_0\|) > 0$  such that for all  $M \geq M_0 := \|\alpha_0\|_\infty$ , for all admissible  $\mathcal{F} \subset \{\alpha : \mathcal{T} \times \mathcal{X} \rightarrow [-M, M]\}$  for any  $\nu, \nu_c > 0$ ,  $0 < \gamma \leq \frac{1}{2}$  and  $n \geq \frac{24M(B_0 \exp M + 1)^2}{\gamma^2(\nu + \nu_c)}$  holds:

$$\begin{aligned} P \left( \sup_{\alpha \in \mathcal{F}} \frac{\Lambda(\alpha_0) - \Lambda(\alpha) - \frac{1}{n}(L_n(\alpha_0) - L_n(\alpha))}{\nu + \nu_c + \Lambda(\alpha_0) - \Lambda(\alpha)} \geq \gamma \right) \\ \leq \kappa(\nu_c, \mathcal{F}) \exp \left( -\frac{\gamma^2\nu n}{\kappa_0 B_0^2 M \exp(2M)} \right), \text{ where} \end{aligned} \quad (2.2)$$

$$\begin{aligned} \kappa(\nu_c, \mathcal{F}) = 6 \sup_{\bar{z} \in \mathcal{Z}^{2n}} \left[ N \left( \frac{\gamma\nu_c}{64 \exp M}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})} \right) N \left( \frac{\gamma\nu_c}{64 \exp(2M)}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])} \right) \right] \\ \text{and } \gamma_0 = \frac{2608}{3} \end{aligned}$$

**Proof:** W.l.g. let  $M_0, B_0 \geq 1$ . Define  $F = \{f_\alpha = \ell(\cdot, \alpha_0) - \ell(\cdot, \alpha); \alpha \in \mathcal{F}\}$ , then by (1.4)  $|f_\alpha| \leq 2(M + \exp M)$ . Also by (1.4)  $Ef_\alpha \geq 0$  and by application of (1.9) and (1.8)

$$\begin{aligned} Ef_\alpha^2 &= E[\ell(\cdot, \alpha_0) - \ell(\cdot, \alpha)]^2 \\ &\leq (B_0 \exp M + 1)^2 \|\alpha - \alpha_0\|_2^2 \\ &\leq 4M(B_0 \exp M + 1)^2 Ef_\alpha. \end{aligned}$$

This implies that the conditions of Theorem 2.1 are fulfilled with  $K_1 = 2(M + \exp M)$ ,  $K_2 = 4M(B_0 \exp M + 1)^2$ .

Therefore, for  $n \geq \max \left\{ 4 \frac{K_1 + K_2}{\gamma^2(\nu + \nu_c)}, \frac{K_1^2}{\gamma^2(\nu + \nu_c)} \right\} = 4 \frac{K_1 + K_2}{\gamma^2(\nu + \nu_c)}$  holds

$$\begin{aligned} P \left( \sup_{\alpha \in \mathcal{F}} \frac{\Lambda(\alpha_0) - \Lambda(\alpha) - \frac{1}{n}(L_n(\alpha_0) - L_n(\alpha))}{\nu + \nu_c + \Lambda(\alpha_0) - \Lambda(\alpha)} \geq \gamma \right) \\ \leq \sup_{\bar{z} \in \mathcal{Z}^{2n}} 2N \left( \frac{\gamma \nu_c}{4}, F, d_{L^1(\nu_{\bar{z}})} \right) \exp \left( -\frac{3\gamma^2 \nu n}{4K_1 + 162K_2} \right) \\ + \sup_{\bar{z} \in \mathcal{Z}^{2n}} 4N \left( \frac{\gamma \nu_c}{4K_1}, F, d_{L^1(\nu_{\bar{z}})} \right) \exp \left( -\frac{\gamma^2 \nu n}{2K_1^2} \right). \end{aligned} \quad (2.3)$$

By easy calculations  $\max \left\{ 2K_1^2, \frac{4K_1 + 162K_2}{3} \right\} \leq \kappa_0 B_0^2 M \exp(2M)$  and  $4 \frac{K_1 + K_2}{\gamma^2(\nu + \nu_c)} \leq n_0 := \frac{24M(B_0 \exp M + 1)^2}{\gamma^2(\nu + \nu_c)}$ . Therefore, using  $4 \exp M \geq K_1 \geq 1$  we obtain that for  $n \geq n_0$  the right hand side of (2.3) is bounded above by

$$\begin{aligned} 6 \sup_{\bar{z} \in \mathcal{Z}^{2n}} N \left( \frac{\gamma \nu_c}{4K_1}, F, d_{L^1(\nu_{\bar{z}})} \right) \exp \left( -\frac{\gamma^2 \nu n}{\kappa_0 B_0^2 M \exp(2M)} \right) \\ \leq 6 \sup_{\bar{z} \in \mathcal{Z}^{2n}} N \left( \frac{\gamma \nu_c}{16 \exp M}, F, d_{L^1(\nu_{\bar{z}})} \right) \exp \left( -\frac{\gamma^2 \nu n}{\kappa_0 B_0^2 M \exp(2M)} \right). \end{aligned}$$

Theorem 2.2 now will be a consequence of the following estimate. For  $\varepsilon > 0$  holds

$$N(\varepsilon, F, d_{L^1(\nu_{\bar{z}})}) \leq N \left( \frac{\varepsilon}{4}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})} \right) N \left( \frac{\varepsilon}{4 \exp M}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])} \right). \quad (2.4)$$

For the proof of (2.4) introduce  $\tilde{F} = \{\tilde{f}_\alpha(\cdot) = \ell(\cdot, \alpha); \alpha \in \mathcal{F}\}$ , then

$$\begin{aligned} N(\varepsilon, F, d_{L^1(\nu_{\bar{z}})}) &\leq N \left( \frac{\varepsilon}{2}, \tilde{F}, d_{L^1(\nu_{\bar{z}})} \right) N \left( \frac{\varepsilon}{2}, \{\ell(\alpha_0)\}, d_{L^1(\nu_{\bar{z}})} \right) \\ &= N \left( \frac{\varepsilon}{2}, \tilde{F}, d_{L^1(\nu_{\bar{z}})} \right). \end{aligned}$$

Define  $\mathcal{H} = \{g_\alpha(x, y, \delta) = \delta \alpha(y, x); \alpha \in \mathcal{F}\}$  and  $\mathcal{K} = \{k_\alpha(x, y, \delta) = \int_0^y \exp \alpha(u, x) du; \alpha \in \mathcal{F}\}$ , then obviously  $N(\varepsilon, \mathcal{H}, d_{L^1(\nu_{\bar{z}})}) \leq N(\varepsilon, \mathcal{F}, d_{L^1(\nu_{\bar{z}})})$ .

Further,

$$\begin{aligned}
d_{L^1(\nu_{\bar{z}})}(k_{\alpha_1}, k_{\alpha_2}) &= \frac{1}{n} \sum_{i=1}^n \int_0^{y_i} |\exp \alpha_1(u, x_i) - \exp \alpha_2(u, x_i)| du \\
&\leq \frac{1}{n} \sum_{i=1}^n \int_0^1 |\exp \alpha_1(u, x_i) - \exp \alpha_2(u, x_i)| du \\
&= d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}(\exp \circ \alpha_1, \exp \circ \alpha_2).
\end{aligned}$$

This implies

$$\begin{aligned}
N(\varepsilon, \mathcal{K}, d_{L^1(\nu_{\bar{z}})}) &\leq N(\varepsilon, \exp \circ \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}) \\
&\leq N\left(\frac{\varepsilon}{\exp M}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}\right)
\end{aligned}$$

using that for  $\mathcal{F}$  with  $|f| \leq K$  for  $f \in \mathcal{F}$  and Lipschitzfunctions  $\varphi : [-K, K] \rightarrow \mathbb{R}$

$$N(\varepsilon, \varphi \circ \mathcal{F}, d_{L^p(\mu)}) \leq N\left(\frac{\varepsilon}{\text{Lip } \varphi}, \mathcal{F}, d_{L^p(\mu)}\right). \quad (2.5)$$

This implies using a well-known upper bound for the covering number of the sum of two function classes

$$\begin{aligned}
N(\varepsilon, \tilde{F}, d_{L^1(\nu_{\bar{z}})}) &= N(\varepsilon, \mathcal{H} \ominus \mathcal{K}, d_{L^1(\nu_{\bar{z}})}) \\
&\leq N\left(\frac{\varepsilon}{2}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})}\right) N\left(\frac{\varepsilon}{2 \exp M}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}\right)
\end{aligned}$$

and we obtain by combining the above estimates the statement of Theorem 2.2.  $\square$

### Remark 2.3 (more general loss functions and estimation problems)

From the proof of Theorem 2.2 one obtains a similar maximal inequality for more general loss functions  $\ell$  (i.e. for more general estimation problems and (minimum contrast) estimation methods) satisfying the following three conditions:

$$\begin{aligned}
|\ell(\alpha_0) - \ell(\alpha)| &\leq K_1 \\
E\ell(\alpha_0) &\geq E\ell(\alpha) \\
E(\ell(\alpha_0) - \ell(\alpha))^2 &\leq K_2 E(\ell(\alpha_0) - \ell(\alpha)).
\end{aligned} \quad (2.6)$$

For (2.6) the following two conditions corresponding to (1.8) and (1.9) are sufficient:

$$E(\ell(\alpha_0) - \ell(\alpha)) \geq k \|\alpha - \alpha_0\|_2^2 \quad (2.7)$$

$$E(\ell(\alpha_0) - \ell(\alpha))^2 \leq \tilde{k} \|\alpha - \alpha_0\|_2^2. \quad (2.8)$$

Therefore, under condition (2.6) we obtain exponential inequalities with

$$\begin{aligned} & N\left(\frac{\gamma\nu_c}{4K_1}, F, d_{L^1(\nu_{\bar{z}})}\right) \quad \text{replacing the capacity term} \\ & N\left(\frac{\gamma\nu_c}{64\exp M}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})}\right) N\left(\frac{\gamma\nu_c}{64\exp(2M)}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}\right) \end{aligned} \quad (2.9)$$

in (2.2), where  $F = \{f_\alpha = \ell(\alpha_0) - \ell(\alpha); \alpha \in \mathcal{F}\}$  is defined as in the proof of Theorem 2.2. This exponential inequality can be applied to prove convergence rates for the corresponding empirical minimum risk estimators. Condition (2.8) corresponds roughly to condition M1 in Birgé and Massart (1998). Condition (2.7) together with an upper bound as in (1.8) corresponds to condition C in Birgé and Massart (1998). Their growth condition M2 involving also the  $L_\infty$ -metric is replaced in our approach by corresponding growth conditions on the  $L^1$ -covering numbers  $N(\cdot, F, d_{L^1(\nu_{\bar{z}})})$  which then is closer related to the  $L^1$ -metric condition with bracketing  $M_{1,[]}$  in Barron, Birgé, and Massart (1999).

### 3 Error bounds for maximum-likelihood estimators for conditional log-hazard functions

As a measure of complexity of a model  $\mathcal{F}$  we define

$$\mathcal{C}_n(\mathcal{F}) = 6 \sup_{\bar{z} \in Z^{2n}} N\left(\frac{1}{n}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})}\right) N\left(\frac{1}{n}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}\right) \quad (3.1)$$

which arises from the first part of the estimate in (2.2). The following theorem estimates the mean  $\Lambda$ -error and the MISE of the ML-estimator in a model  $\mathcal{F}$ . Admissibility of  $\mathcal{F}$  is a weak measurability condition (cf. Lee, Bartlett, and Williamson (1996)) which is satisfied for the examples considered in this paper.

**Theorem 3.1** *Let  $\mathcal{F} \subset \{\alpha : \mathcal{T} \times \mathcal{X} \rightarrow [-M, M]\}$  be admissible where  $M \geq M_0 = \|\alpha_0\|_\infty$  and  $B_0, \kappa_0$  are as in Theorem 2.2. Assume that  $\mathcal{C}_n(\mathcal{F}) < \infty$ , then for the ML-estimator  $\hat{\alpha}_n = \arg \max_{\alpha \in \mathcal{F}} L_n(\alpha)$  the following error estimates hold:*

$$\begin{aligned} & E |\Lambda(\hat{\alpha}_n) - \Lambda(\alpha_0)| \\ & \leq 2 \inf_{\alpha \in \mathcal{F}} |\Lambda(\alpha) - \Lambda(\alpha_0)| + 8\kappa_0 B_0^2 M \exp(2M) \frac{\log \mathcal{C}_n(\mathcal{F}) + 1}{n} \end{aligned} \quad (3.2)$$

and

$$\begin{aligned} & E \|\hat{\alpha}_n - \alpha_0\|_2^2 \\ & \leq 2 \exp(2M) \inf_{\alpha \in \mathcal{F}} \|\alpha - \alpha_0\|_2^2 + 32\kappa_0 B_0^2 M^2 \exp(2M) \frac{\log \mathcal{C}_n(\mathcal{F}) + 1}{n}. \end{aligned} \quad (3.3)$$



**Proof:** In our proof we use a similar technique as in the context of regression estimation in Kohler(1997, 1999a, 1999 b). We decompose the  $\Lambda$ -error into two parts

$$|\Lambda(\widehat{\alpha}_n) - \Lambda(\alpha_0)| = T_{1,n} + T_{2,n} \quad (3.4)$$

with  $T_{1,n} = \Lambda(\alpha_0) - \Lambda(\widehat{\alpha}_n) - \frac{2}{n}(L_n(\alpha_0) - L_n(\widehat{\alpha}_n))$  and  $T_{2,n} = \frac{2}{n}(L_n(\alpha_0) - L_n(\widehat{\alpha}_n))$ . From the definition of  $\widehat{\alpha}_n$  we obtain by a standard argument  $ET_{2,n} \leq 2 \inf_{\alpha \in \mathcal{F}} |\Lambda(\alpha) - \Lambda(\alpha_0)|$ .

It remains to establish

$$ET_{1,n} \leq 8\kappa_0 B_0^2 M \exp(2M) \frac{\log \mathcal{C}_n(\mathcal{F}) + 1}{n}. \quad (3.5)$$

For  $t \geq t_0 = \frac{96M(B_0 \exp M + 1)^2}{n}$  we obtain from Theorem 2.2 with  $\gamma = \frac{1}{2}$ ,  $\nu = \nu_c = \frac{t}{2}$

$$\begin{aligned} P(T_{1,n} \geq t) &\leq P\left(\sup_{\alpha \in \mathcal{F}} \frac{\Lambda(\alpha_0) - \Lambda(\alpha) - \frac{1}{n}(L_n(\alpha_0) - L_n(\alpha))}{\frac{t}{2} + \frac{t}{2} + \Lambda(\alpha_0) - \Lambda(\alpha)} \geq \frac{1}{2}\right) \\ &\leq 6 \sup_{\bar{z} \in \mathcal{Z}^{2n}} \left[ N\left(\frac{\frac{1}{2} \frac{t_0}{2}}{64 \exp M}, \mathcal{F}, d_{L^1(\nu_{\bar{z}})}\right) N\left(\frac{\frac{1}{2} \frac{t_0}{2}}{64 \exp(2M)}, \mathcal{F}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])}\right) \right] \\ &\quad \cdot \exp\left(-\frac{tn}{8\kappa_0 B_0^2 M \exp(2M)}\right). \end{aligned}$$

For  $M \geq 1$  holds  $\frac{\frac{1}{2} \frac{t_0}{2}}{64 \exp(2M)} \geq \frac{1}{n}$  and, therefore,

$$P(T_{1,n} \geq t) \leq \mathcal{C}_n(\mathcal{F}) \exp\left(-\frac{tn}{8\kappa_0 B_0^2 M \exp(2M)}\right).$$

This implies for  $u \geq t_0$

$$\begin{aligned} ET_{1,n} &\leq \int_0^u 1 \, dt + \int_u^\infty P(T_{1,n} \geq t) \, dt \\ &\leq u + \mathcal{C}_n(\mathcal{F}) \frac{8\kappa_0 B_0^2 M \exp(2M)}{n} \exp\left(-\frac{un}{8\kappa_0 B_0^2 M \exp(2M)}\right). \end{aligned} \quad (3.6)$$

The r.h.s. of (3.6) is minimized by  $u_0 = \frac{1}{n} 8\kappa_0 B_0^2 M \exp(2M) \log \mathcal{C}_n(\mathcal{F})$ . It is easy to see that  $u_0 \geq t_0$ . With this  $u_0$  inserted in (3.6) we obtain inequality (3.5) and so statement (3.2).

From (1.8) we then conclude

$$\begin{aligned} E \|\widehat{\alpha}_n - \alpha_0\|_2^2 &\leq 4M E |\Lambda(\widehat{\alpha}_n) - \Lambda(\alpha_0)| \\ &\leq 8M \inf_{\alpha \in \mathcal{F}} |\Lambda(\alpha) - \Lambda(\alpha_0)| + 32\kappa_0 B_0^2 M^2 \exp(2M) \frac{\log \mathcal{C}_n(\mathcal{F}) + 1}{n} \\ &\leq 8M \frac{\exp(2M)}{4M^2} \inf_{\alpha \in \mathcal{F}} \|\alpha - \alpha_0\|_2^2 + 32\kappa_0 B_0^2 M^2 \exp(2M) \frac{\log \mathcal{C}_n(\mathcal{F}) + 1}{n}. \end{aligned} \quad \square$$

**Remark 3.2 (general estimation problem)**

The error estimates in Theorem 3.1 decompose the error as usual in an approximation error and a stochastic error of the order  $\frac{\log \mathcal{C}_n(\mathcal{F})+1}{n}$ . As in Remark 2.3 (see (2.9)) we obtain a similar estimate for general loss functions  $\ell$  by replacing the model complexity term  $\mathcal{C}_n(\mathcal{F})$  by

$$\mathcal{C}_n(F) = 6 \sup_{\bar{z} \in Z^{2n}} N\left(\frac{1}{n}, F, d_{L^1(\nu_{\bar{z}})}\right) \quad (3.7)$$

with  $F = \{f_\alpha = \ell(\cdot, \alpha_0) - \ell(\cdot, \alpha); \alpha \in \mathcal{F}\}$ . In comparison to a related result in Birgé and Massart (1998, Corollary 1, section 5) which uses in condition M2 assumptions on the  $L^2$ - and  $L^\infty$ -covering numbers of  $\mathcal{F}$  our estimate uses only  $L^1$ -covering numbers in the model complexity term  $\mathcal{C}_n(\mathcal{F})$  resp.  $\mathcal{C}_n(F)$ . Our condition is closer to the  $L^1$ -condition with bracketing  $M_{1,[]}$  in Barron, Birgé, and Massart (1999, section 6).

From Pollard's estimate for bounded VC-classes  $\mathcal{F}$ ,  $d = \dim_{VC} \mathcal{F}$ , with majorant  $H$  stating that for  $\varepsilon > 0$

$$N\left(\varepsilon \|H\|_{L^p(\mu)}, \mathcal{F}, d_{L^p(\mu)}\right) \leq \kappa d (16e)^d \left(\frac{1}{\varepsilon}\right)^{p(d-1)}, \quad (3.8)$$

(see van der Vaart and Wellner (1996, Theorem 2.6.7)) we obtain from our estimate in (3.3) a direct connection of convergence rates to the VC-dimension of the class  $\mathcal{F}$ .

As a consequence of this remark we obtain

**Corollary 3.3** *Under the conditions of Theorem 3.1 where  $\mathcal{F}$  is a bounded VC-class we obtain*

$$E \|\hat{\alpha}_n - \alpha_0\|_2^2 \leq C_1(M) \inf_{\alpha \in \mathcal{F}} \|\alpha - \alpha_0\|_2^2 + C_2(M, B_0) \dim_{VC}(\mathcal{F}) \frac{\log n}{n}. \quad (3.9)$$

A similar convergence rate result holds for general estimation problems as in Remarks 2.3, 3.2.

**Remark 3.4 (Sieve estimators)** *Let  $(\mathcal{F}_K)_{K \in \mathbb{N}}$  be a sieve of VC-classes in the underlying model  $\mathcal{F}$  with  $D_K = \dim_{VC} \mathcal{F}_K$  and approximation rate  $b_K = \inf_{\alpha \in \mathcal{F}_K} \|\alpha - \alpha_0\|_2^2$ . Assume that for some  $r, s > 0$*

$$b_K = O(K^{-r}), D_K = O(K^s). \quad (3.10)$$

*There are two well studied types of sieves, linear sieves, i.e. finite dimensional vector spaces which approximate typically smooth function classes and secondly nets (like neural nets, radial basis function nets, ...). Under assumption (3.10) we obtain from the estimate in (3.9) when choosing the optimal*

parameter  $K_n$  in the bias-variance decomposition (3.9) an estimate for the MISE of  $\hat{\alpha}_n$  of the form:

$$E\|\hat{\alpha}_n - \alpha_0\|_2^2 = O\left(\left(\frac{\log n}{n}\right)^{\frac{r}{r+s}}\right). \quad (3.11)$$

Here  $r$  determines the approximation rate of the sieve which is usually for splines, wavelets, polynomials related to smoothness of the parameter and  $s$  determines the complexity of the net.

If  $\mathcal{F}_K$  is a subset of a  $K$ -dimensional vector space, then  $s = 1$ , and if  $r = \frac{2p}{d}$  (with  $p = \text{degree of smoothness}$ ,  $d = \text{dimension of space}$ ) we will obtain in some examples optimal convergence rates up to logarithmic terms.

Polynomial rates (i.e.  $b_K$  are as in (3.10)) can also be obtained for function classes which are obtained from VC-classes by some operations like transformations, sums, etc.

## 4 Structural risk-minimization

Based on the maximal inequality in Theorem 2.2 one obtains for  $\eta \in (0, 1)$  and any data dependent estimator  $\alpha_n \in \mathcal{F}$ , that with probability  $1 - \eta$

$$\begin{aligned} & |\Lambda(\alpha_n) - \Lambda(\alpha_0)| \\ & \leq \frac{1}{n} 8\kappa_0 B_0^2 M \exp(2M) \log \frac{\mathcal{C}_n(\mathcal{F})}{\eta} + \frac{2}{n} (L_n(\alpha_0) - L_n(\alpha_n)). \end{aligned} \quad (4.1)$$

The idea of structural risk minimization (complexity regularization) due to Vapnik (1995) is to construct an estimator minimizing approximatively the r.h.s. of (4.1), i.e. minimizing

$$c_n \frac{\log \mathcal{C}_n(\mathcal{F})}{n} - \frac{2}{n} L_n(\alpha) \quad (4.2)$$

where  $c_n$  is a slowly increasing function independent of the unknown parameters which asymptotically majorizes the corresponding constant in (4.1). The minimization is carried out not only over  $\alpha$  in one fixed class  $\mathcal{F} = \mathcal{F}_n$  but allows to choose  $\alpha$  within a finite set of model classes  $\{\mathcal{F}_{n,p}; p \in \mathcal{P}_n\}$ ,  $p$  typically describing some smoothness or network complexity. The error term  $c_n \frac{\log \mathcal{C}_n(\mathcal{F}_{n,p})}{n}$  can be interpreted as a penalization term for the complexity of the model.

A detailed and general description of this approach with several applications has been given in Barron, Birgé, and Massart (1999) based on the error estimates in Birgé and Massart (1998) as well as on new tools. In that paper one also finds several references to this method. In our paper we use some technical ideas from Kohler (1997, proof of Theorem 4.2), concerning regression estimates minimizing empirical penalized squared loss there.

Let  $M_0 = \|\alpha_0\|_\infty > 0$ ,  $B_0 = B_0(\|\alpha_0\|_\infty) > 0$  be as in Theorem 2.2, and let  $\mathcal{P}_n$  be finite sets for  $n \in \mathbb{N}$ , and for  $p \in \mathcal{P}_n$  let  $\mathcal{F}_{n,p} \subset \{\alpha : \mathcal{T} \times \mathcal{X} \rightarrow [-M, M]\}$  be admissible models,  $M \geq M_0$  with  $\mathcal{C}_n(\mathcal{F}_{n,p}) < \infty$ ,  $\forall p \in \mathcal{P}_n$ . Then the complexity regularized estimator  $\alpha_n^*$  is defined in two steps:

$$1. \quad \text{Let } p_n^* = \arg \min_{p \in \mathcal{P}_n} \left( -\frac{1}{n} \sup_{\alpha \in \mathcal{F}_{n,p}} L_n(\alpha) + \text{pen}_n(p) \right) \quad (4.3)$$

where  $\text{pen}_n(p)$  is a penalization term for the complexity of model  $\mathcal{F}_{n,p}$  satisfying asymptotically as  $n \rightarrow \infty$

$$\text{pen}_n(p) \geq 4\kappa_0 B_0^2 M \exp(2M) \frac{\log \mathcal{C}_n(\mathcal{F}_{n,p})}{n}. \quad (4.4)$$

$$2. \quad \alpha_n^* = \arg \max_{\alpha \in \mathcal{F}_{n,p_n^*}} L_n(\alpha). \quad (4.5)$$

It is important to note that the r.h.s. of (4.4) is not supposed to be the actual penalty term used in application since it depends on the unknown  $M_0$  and  $B_0$ . This expression represents a lower bound for the penalty, sufficient for Theorem 4.1 to hold(cf. also (4.2)). For asymptotic results the actual penalty term should be chosen independently of  $M_0$  and  $B_0$ , majorising the r.h.s. of (4.4) for large sample sizes. An example of how this can be done is given in Theorem 5.3. The following theorem gives an error bound for complexity regularized sieve estimators based on the maximal inequality in Theorem 2.2. A general related error bound is given in Barron, Birgé, and Massart (1999, Theorem 8) under some alternative conditions on  $L_2 - L_\infty$  covering respectively  $L_1$  covering with bracketing.

**Theorem 4.1** *For the complexity regularized ML-estimator  $\alpha_n^*$  the following error estimates hold:*

$$\begin{aligned} E |\Lambda(\alpha_n^*) - \Lambda(\alpha_0)| &\leq 2 \inf_{p \in \mathcal{P}_n} \left( \text{pen}_n(p) + \inf_{\alpha \in \mathcal{F}_{n,p}} |\Lambda(\alpha) - \Lambda(\alpha_0)| \right) \\ &\quad + \frac{4\kappa_0 B_0^2 M \exp(2M)}{n} (1 + \log |\mathcal{P}_n|) \end{aligned} \quad (4.6)$$

and

$$\begin{aligned} E \|\alpha_n^* - \alpha_0\|_2^2 &\leq 2 \inf_{p \in \mathcal{P}_n} \left( 4M \text{pen}_n(p) + \exp(2M) \inf_{\alpha \in \mathcal{F}_{n,p}} \|\alpha - \alpha_0\|_2^2 \right) \\ &\quad + \frac{16\kappa_0 B_0^2 M^2 \exp(2M)}{n} (1 + \log |\mathcal{P}_n|). \end{aligned} \quad (4.7)$$

**Proof:** As in the proof of Theorem 3.1 we consider the decomposition of the error into two terms

$$\begin{aligned} T_{1,n} &:= \Lambda(\alpha_0) - \Lambda(\alpha_n^*) - \frac{2}{n} (L_n(\alpha_0) - L_n(\alpha_n^*)) - 2 \text{pen}_n(p_n^*) \\ T_{2,n} &:= \frac{2}{n} (L_n(\alpha_0) - L_n(\alpha_n^*)) + 2 \text{pen}_n(p_n^*). \end{aligned} \quad (4.8)$$

Our first aim is to prove

$$ET_{1,n} \leq \frac{4\kappa_0 B_0^2 M \exp(2M)}{n} (1 + \log |\mathcal{P}_n|). \quad (4.9)$$

For the proof we obtain as in Kohler (1997, pg. 85).

$$P(T_{1,n} > t) \leq \sum_{p \in \mathcal{P}_n} P \left( \sup_{\alpha \in \mathcal{F}_{n,p}} \frac{\Lambda(\alpha_0) - \Lambda(\alpha) - \frac{1}{n}(L_n(\alpha_0) - L_n(\alpha))}{t + 2 \text{pen}_n(p) + \Lambda(\alpha_0) - \Lambda(\alpha)} \geq \frac{1}{2} \right).$$

Then for  $n \geq \frac{24M(B_0 \exp M + 1)^2}{\frac{1}{4}t}$  we obtain from Theorem 2.2 with  $\gamma = \frac{1}{2}$ ,  $\nu = t + \text{pen}_n(p)$ ,  $\nu_c = \text{pen}_n(p)$  observing that the condition  $n \geq \frac{24M(B_0 \exp M + 1)^2}{\frac{1}{4}(t + 2 \text{pen}_n(p))}$  is fulfilled for any  $p \in \mathcal{P}_n$

$$\begin{aligned} P(T_{1,n} > t) &\leq \sum_{p \in \mathcal{P}_n} \underbrace{\left[ 6 \sup_{\bar{z} \in \mathcal{Z}^{2n}} N \left( \frac{\frac{1}{2} \text{pen}_n(p)}{64 \exp M}, \mathcal{F}_{n,p}, d_{L^1(\nu_{\bar{z}})} \right) N \left( \frac{\frac{1}{2} \text{pen}_n(p)}{64 \exp(2M)}, \mathcal{F}_{n,p}, d_{L^1(\tilde{\nu}_{\bar{z}} \otimes U[0,1])} \right) \right]}_{=: s_n(p)} \\ &\quad \cdot \exp \left( - \frac{\text{pen}_n(p)n}{4\kappa_0 B_0^2 M \exp(2M)} \right) \exp \left( - \frac{tn}{4\kappa_0 B_0^2 M \exp(2M)} \right). \end{aligned}$$

Since  $\log \mathcal{C}_n(\mathcal{F}_{n,p}) \geq 1$  and hence  $\frac{1}{n} \leq \frac{\frac{1}{2} \text{pen}_n(p)}{64 \exp M}, \frac{\frac{1}{2} \text{pen}_n(p)}{64 \exp(2M)}$  for any  $p \in \mathcal{P}_n$ , we obtain  $s_n(p) \leq \mathcal{C}_n(\mathcal{F}_{n,p})$ . Further

$$\begin{aligned} s_n(p) \exp \left( - \frac{\text{pen}_n(p)}{4\kappa_0 B_0^2 M \exp(2M)} n \right) &\leq \exp \left( \log \mathcal{C}_n(\mathcal{F}_{n,p}) - \frac{\text{pen}_n(p)}{4\kappa_0 B_0^2 M \exp(2M)} n \right) \\ &\leq \exp(0) = 1 \end{aligned}$$

by definition of  $\text{pen}_n(p)$ , and, therefore, for  $t \geq t_0 := \frac{96M(B_0 \exp M + 1)^2}{n}$  holds

$$P(T_{1,n} > t) \leq |\mathcal{P}_n| \exp \left( - \frac{t}{4\kappa_0 B_0^2 M \exp(2M)} n \right). \quad (4.10)$$

This implies for  $u \geq t_0$

$$\begin{aligned} ET_{1,n} &\leq \int_0^u 1 \, dt + \int_u^\infty P(T_{1,n} \geq t) dt \\ &\leq u + |\mathcal{P}_n| \frac{4\kappa_0 B_0^2 M \exp(2M)}{n} \exp \left( - \frac{u}{4\kappa_0 B_0^2 M \exp(2M)} n \right). \end{aligned} \quad (4.11)$$

The r.h.s. of this inequality is minimized by

$$u = u_0 := \frac{4\kappa_0 B_0^2 M \exp(2M)}{n} \log |\mathcal{P}_n|,$$

and w.l.g. for  $\kappa_0 \cdot \log |\mathcal{P}_n| \geq 24(1 + \frac{1}{B_0 \exp M})^2$  it holds that  $u_0 \geq t_0$ . This choice of  $u$  leads to (4.9).

From the definition of  $p_n^*$  and  $\alpha_n^*$  we obtain that

$$\begin{aligned} T_{2,n} &= 2 \left[ \frac{1}{n} L_n(\alpha_0) - \frac{1}{n} \sup_{\alpha \in \mathcal{F}_{n,p_n^*}} L_n(\alpha) + \text{pen}_n(p_n^*) \right] \\ &= 2 \left[ \frac{1}{n} L_n(\alpha_0) + \inf_{p \in \mathcal{P}_n} \left( -\frac{1}{n} \sup_{\alpha \in \mathcal{F}_{n,p}} L_n(\alpha) + \text{pen}_n(p) \right) \right] \\ &= 2 \inf_{p \in \mathcal{P}_n} \left[ \inf_{\alpha \in \mathcal{F}_{n,p}} \frac{1}{n} (L_n(\alpha_0) - L_n(\alpha)) + \text{pen}_n(p) \right]. \end{aligned} \quad (4.12)$$

Therefore, using that  $\text{pen}_n(p)$  is deterministic we obtain

$$\begin{aligned} ET_{2,n} &\leq 2 \inf_{p \in \mathcal{P}_n} E \left[ \inf_{\alpha \in \mathcal{F}_{n,p}} \frac{1}{n} (L_n(\alpha_0) - L_n(\alpha)) + \text{pen}_n(p) \right] \\ &\leq 2 \inf_{p \in \mathcal{P}_n} \left[ \inf_{\alpha \in \mathcal{F}_{n,p}} E \frac{1}{n} (L_n(\alpha_0) - L_n(\alpha)) + \text{pen}_n(p) \right] \\ &\leq 2 \inf_{p \in \mathcal{P}_n} \left[ \inf_{\alpha \in \mathcal{F}_{n,p}} |\Lambda(\alpha) - \Lambda(\alpha_0)| + \text{pen}_n(p) \right]. \end{aligned} \quad (4.13)$$

(4.11) and (4.13) imply (4.6). The proof of (4.7) then follows from (1.8).  $\square$

## 5 Adaptive sieve estimates for the conditional log-hazard function

In this section we apply the results of sections 3 and 4 to several types of sieves. In the first part we obtain that the complexity regularized spline estimate is approximatively optimal even with unknown degree of smoothness, i.e. it has up to a logarithmic term the same optimal convergence rate as the estimator of Kooperberg, Stone, and Truong (1995b) in the case with known degree of smoothness. In the second part we obtain convergence results for net estimates under the assumption that the conditional log-hazards have a certain representation property. Some applications to Sobolev class models will be considered in a subsequent paper.

### 5.1 Tensor product splines

In this section we consider tensor product splines. For general background of this class of functions we refer to Kohler (1997) and the references given there.

Let  $V_{h,M}$  denote the class of tensor product splines of  $[-hM, 1 + hM]^{k+1}$  of degree  $M \in \mathbb{N}_0$  in each coordinate and of grid width  $h > 0$ . Let  $\Phi(L, V_{h,M})$  (for  $L > 0$ ) denote the class of truncated functions  $T_L \circ g$ ,  $g \in V_{h,M}$  where

$$T_L \circ g = \begin{cases} L & \text{if } g \geq L \\ g & \text{if } -L \leq g \leq L \\ -L & \text{if } g \leq -L. \end{cases} \quad (5.1)$$

We consider for  $p = r + \beta$ ,  $r \in \mathbb{N}_0$ ,  $\beta \in (0, 1)$  the smoothness classes  $\Sigma(p, L)$  of bounded conditional hazard functions  $\alpha(t, x)$  on  $[0, 1]^{k+1}$  satisfying for all  $z_1, z_2 \in [0, 1]^{k+1}$  a Hölder condition of order  $p$

$$\|\alpha\|_\infty \leq L \quad \text{and} \quad |D^r \alpha(z_1) - D^r \alpha(z_2)| \leq L \|z_1 - z_2\|_2^\beta. \quad (5.2)$$

For classes with known degree of smoothness we obtain the following result.

**Theorem 5.1 (known smoothness class)** *Let  $1 \leq p < \infty$ ,  $L > 0$ ,  $\widetilde{M} \in \mathbb{N}$ ,  $\widetilde{M} \geq p - 1$  and  $h_n = \left(\frac{\log n}{n}\right)^{\frac{1}{2p+k+1}}$ . Then the spline ML-estimator*

$$\widehat{\alpha}_n = \arg \max_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} L_n(\alpha) \quad (5.3)$$

*satisfies*

$$\sup_{\alpha \in \Sigma(p, L)} E \|\widehat{\alpha}_n - \alpha_0\|_2^2 = O \left( \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right) \quad (5.4)$$

*and*

$$\sup_{\alpha \in \Sigma(p, L)} E |\Lambda(\widehat{\alpha}_n) - \Lambda(\alpha_0)| = O \left( \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right). \quad (5.5)$$

**Proof:** From definition of the truncation operator  $T_L$  it follows that

$$\inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} \|\alpha_0 - \alpha\|_\infty \leq \inf_{\alpha \in V_{h_n, \widetilde{M}}} \|\alpha_0 - \alpha\|_\infty$$

and, therefore, from the approximation result in Kohler (1997, Lemma 1.5) for approximation of Hölder-continuous functions by tensor product splines (which is in sup-norm) we obtain

$$\inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} \|\alpha_0 - \alpha\|_2^2 \leq C h_n^{2p} \leq C \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \quad (5.6)$$

with  $C = C(p, L)$  independent of  $\alpha_0$ .

To estimate the stochastic error in Theorem 3.1 note that  $V_{h_n, \widetilde{M}}$  is a vector space of dimension  $\leq \left(\left\lceil \frac{1}{h_n} \right\rceil + \widetilde{M}\right)^{k+1}$  (see Kohler (1997, pg. 79) and, therefore, (cf. van der Vaart and Wellner (1996, Lemma 2.6.18))

$$\dim_{VC} \Phi(L, V_{h_n, \widetilde{M}}) \leq \dim_{VC} V_{h_n, \widetilde{M}} \leq \left(\left\lceil \frac{1}{h_n} \right\rceil + \widetilde{M}\right)^{k+1} + 2. \quad (5.7)$$

Therefore, we obtain from Theorem 3.1

$$E \|\widehat{a}_n - a_0\|_2^2 \leq C_1 \left(\frac{\log n}{n}\right)^{\frac{2p}{2p+k+1}} + C_2 \left[ \left(\left\lceil \left(\frac{n}{\log n}\right)^{\frac{1}{2p+k+1}} \right\rceil + \widetilde{M}\right)^{k+1} + 2 \right] \frac{\log n}{n}. \quad (5.8)$$

This implies with  $M := L$ ,  $\mathcal{F} := \Phi(L, V_{h_n, \widetilde{M}})$  and  $B_0 = B_0(L)$  the first statement (5.4).

For the proof of (5.5) we next establish the approximation rate

$$\inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} |\Lambda(\alpha) - \Lambda(\alpha_0)| \leq C h_n^{2p} \quad (5.9)$$

for the  $\Lambda$ -distance.

From the representation (1.4) and using some elementary properties of the function  $G$  we obtain

$$\begin{aligned} \inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} |\Lambda(\alpha) - \Lambda(\alpha_0)| &= \inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} \int_{T \times \mathcal{X}} \overline{F}_{C|X} G(\alpha - \alpha_0) dP^{(T, X)} \\ &\leq \inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} \int_{T \times \mathcal{X}} \overline{F}_{C|X} G(\|\alpha - \alpha_0\|_\infty) dP^{(T, X)} \\ &\leq \inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} G(\|\alpha - \alpha_0\|_\infty) \\ &\leq G\left(\inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} \|\alpha - \alpha_0\|_\infty\right) \end{aligned}$$

For the last inequality we observe that  $\inf_{x \in A} G(x) = G(\inf A)$  for  $A \subset \mathbb{R}_+$ . Since for  $x_n \downarrow 0$ ,  $G(x_n) = O(x_n^2)$ , we obtain  $\inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} |\Lambda(\alpha) - \Lambda(\alpha_0)| = O((\inf_{\alpha \in \Phi(L, V_{h_n, \widetilde{M}})} \|\alpha - \alpha_0\|_\infty)^2) = O((\inf_{\alpha \in V_{h_n, \widetilde{M}}} \|\alpha_0 - \alpha\|_\infty)^2) = O(h_n^{2p})$  as in (5.6).  $\square$

**Remark 5.2** *The convergence rate in (5.4) for the MISE is up to a logarithmic factor optimal in the minimax-sense (see Kooperberg, Stone, and Truong (1995b, Remark to Corollary 1, note however that the convergence in MISE is stronger), i.e.*

$$\liminf_{n \rightarrow \infty} n^{\frac{2p}{2p+k+1}} \inf_{\widehat{\alpha}_n} \sup_{\alpha_0 \in \Sigma(p, L)} E \|\widehat{\alpha}_n - \alpha_0\|_2^2 > 0 \quad (5.10)$$

for any  $p \geq 1$ ,  $L > 0$ .



If we do not know the smoothness parameter, i.e. assume that

$$\alpha \in \Sigma := \bigcup_{1 \leq p < \infty; L < \infty} \Sigma(p, L) \quad (5.11)$$

then we will obtain that our penalized spline ML-estimator defined in (4.5) adapts up to a logarithmic factor to the unknown smoothness and is up to  $(\log n)^2$  minimax adapted in the sense of Barron, Birgé, and Massart (1999). An *adaptive* estimation method ('HARE') had been introduced in Kooperberg, Stone, and Truong (1995b) and empirically investigated there, however no adaptation result has been proved. As result we obtain that the complexity regularized estimator  $\alpha_n^*$  from (4.5) is approximately adaptive.

**Theorem 5.3 (unknown smoothness degree, adaptation)**

For  $n \in \mathbb{N}$ ,  $q_{\max}(n)$ ,  $K_{\max}(n) \in \mathbb{N}$  let  $\mathcal{P}_n := \{(K, q) \in \mathbb{N} \times \mathbb{N} \mid K \leq K_{\max}(n), q \leq q_{\max}(n)\}$  and for  $(K, q) \in \mathcal{P}_n$  and  $\beta_n := \frac{1}{5} \log \log n$  define the models  $\mathcal{F}_{n,(K,q)} := \Phi(\beta_n, V_{\frac{1}{K}, q-1})$ . Define the complexity regularized estimate  $\alpha_n^*$  as in (4.5) with penalization term

$$\text{pen}_n((K, q)) := \frac{(\log n)^{\frac{8}{5}}}{n} [(K + q - 1)^{k+1} + 2].$$

With  $K_{\max}(n) := n$  and  $q = q_{\max}(n) \rightarrow \infty$  such that  $\frac{q_{\max}(n)}{n} \rightarrow 0$  we obtain for  $p \geq 1$  and  $L > 0$

$$\sup_{\alpha_0 \in \Sigma(p, L)} E \|\alpha_n^* - \alpha_0\|_2^2 = O \left( \log n \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right) \quad (5.12)$$

and

$$\sup_{\alpha_0 \in \Sigma(p, L)} E |\Lambda(\alpha_n^*) - \Lambda(\alpha_0)| = O \left( \log n \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right). \quad (5.13)$$

**Proof:** For the proof of (5.12), (5.13) we first establish the following more general estimates:

For  $p \geq 1$ ,  $L > 0$  there exists  $N_0 = N_0(L) \in \mathbb{N}$  such that for any  $n \geq N_0$ :

$$\begin{aligned} & \sup_{\alpha_0 \in \Sigma(p, L)} E \|\alpha_n^* - \alpha_0\|_2^2 \\ & \leq 2 \inf_{(K, q) \in \mathcal{P}_n} \left( \frac{4 \log \log n (\log n)^{\frac{8}{5}}}{n} [(K + q - 1)^{k+1} + 2] \right. \\ & \quad \left. + (\log n)^{\frac{2}{5}} \inf_{\alpha \in \mathcal{F}_{n,(K,q)}} \|\alpha - \alpha_0\|_2^2 \right) \\ & \quad + \frac{16\kappa_0 B_0^2 (\log \log n)^2 (\log n)^{\frac{2}{5}}}{25 n} (1 + \log(K_{\max}(n) q_{\max}(n))) \end{aligned} \quad (5.14)$$

and

$$\begin{aligned}
& \sup_{\alpha_0 \in \Sigma(p, L)} E |\Lambda(\alpha_n^*) - \Lambda(\alpha_0)| \\
& \leq 2 \inf_{(K, q) \in \mathcal{P}_n} \left( \frac{(\log n)^{\frac{8}{5}}}{n} [(K + q - 1)^{k+1} + 2] \right. \\
& \quad \left. + \inf_{\alpha \in \mathcal{F}_{n, (K, q)}} |\Lambda(\alpha_n^*) - \Lambda(\alpha)| \right) \\
& \quad + \frac{16\kappa_0 B_0^2 \log \log n (\log n)^{\frac{2}{5}}}{5n} (1 + \log(K_{\max}(n) q_{\max}(n))),
\end{aligned} \tag{5.15}$$

where  $\kappa_0$  and  $B_0$  are as in Theorem 3.1.

The statements (5.14), (5.15) follow from Theorem 4.1 with  $M := \beta_n > M_0 := L, n \geq n_0$  and the estimate

$$\text{pen}_n((K, q)) \geq 4\kappa_0 B_0^2 \beta_n \exp(2\beta_n) \frac{\log \mathcal{C}_n(\mathcal{F}_{n, (K, q)})}{n}. \tag{5.16}$$

For the proof of (5.16) we use that for a  $K$ -dimensional vector space  $V$  of functions and  $\beta > 0$  we have for any probability measure  $\mu$  on  $\mathcal{T} \times \mathcal{X}$ ,  $\varepsilon > 0$  the estimate:

$$N(\varepsilon, \Phi(\beta, V), d_{L^p(\mu)}) \leq \kappa(K + 2)(16e)^{K+2} \beta^{p(K+1)} \left( \frac{1}{\varepsilon} \right)^{p(K+1)} \tag{5.17}$$

with some universal constant  $\kappa$ .

This implies that

$$\begin{aligned}
& N(\varepsilon, \mathcal{F}_{n, (K, q)}, d_{L^1(\mu)}) \\
& \leq \kappa((K + q - 1)^{k+1} + 2)(16e)^{(K+q-1)^{k+1}+2} \left( \frac{\beta_n}{\varepsilon} \right)^{(K+q-1)^{k+1}+1},
\end{aligned} \tag{5.18}$$

since  $V_{\frac{1}{K}, q-1}$  has dimension  $\leq (K + q - 1)^{k+1}$ . Therefore

$$\begin{aligned}
\log \mathcal{C}_n(\mathcal{F}_{n, (K, q)}) & \leq \log \kappa^2 + 2 \log((K + q - 1)^{k+1} + 2) \\
& \quad + [(K + q - 1)^{k+1} + 2] [\log(n\beta_n) + 2 \log 16e] \\
& \leq \log \kappa^2 \\
& \quad + [(K + q - 1)^{k+1} + 2] [\log(n\beta_n) + 2 \log 16e + 2] \\
& \leq 2 [(K + q - 1)^{k+1} + 2] \log(n\beta_n)
\end{aligned}$$

for  $n \geq N_0$  where  $N_0$  is independent of  $K, q, L, p$ . This implies

$$\begin{aligned}
\frac{4\kappa_0 B_0^2 \beta_n \exp(2\beta_n) \frac{\log \mathcal{C}_n(\mathcal{F}_{n, (K, q)})}{n}}{\text{pen}_n((K, q))} & \leq \frac{4\kappa_0 B_0^2 \beta_n \exp(2\beta_n) 2 \log(n\beta_n)}{(\log n)^{\frac{8}{5}}} \\
& \leq \frac{4\kappa_0 B_0^2 \beta_n \exp(2\beta_n) 4}{(\log n)^{\frac{3}{5}}} \\
& \leq 1
\end{aligned}$$

for  $n \geq N_0(\kappa_0, B_0(L))$ , and so the result follows.

For the proof of (5.12) let  $K_n := \left\lceil \left( \frac{n}{\log n} \right)^{\frac{1}{2p+k+1}} \right\rceil$ . Then by the approximation result in (5.2) for  $K_n \leq K_{\max}(n)$ , and for  $q_{\max}(n) \geq p$  holds

$$\begin{aligned}
& \inf_{(K,q) \in \mathcal{P}_n} \left( 4\beta_n \text{pen}_n((K, q)) + \exp(2\beta_n) \inf_{\alpha \in \mathcal{F}_{n,(K,q)}} \|\alpha - \alpha_0\|_2^2 \right) \\
& \leq 4\beta_n \text{pen}_n((K_n, p)) + \exp(2\beta_n) \inf_{\alpha \in \mathcal{F}_{n,(K_n,p)}} \|\alpha - \alpha_0\|_2^2 \\
& \leq C_1 \frac{\log \log n (\log n)^{\frac{8}{5}}}{n} [(K_n + p - 1)^{k+1} + 2] + C_2 (\log n)^{\frac{2}{5}} \left( \frac{1}{K_n} \right)^{2p} \\
& = O \left( \frac{\log \log n (\log n)^{\frac{8}{5}}}{n} \left( \frac{n}{\log n} \right)^{\frac{k+1}{2p+k+1}} \right) + O \left( (\log n)^{\frac{2}{5}} \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right) \\
& = O \left( \log n \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right).
\end{aligned}$$

This yields an estimate for the first term in (5.14). For the second term we use the assumptions on  $K_{\max}(n)$  and  $q_{\max}(n)$  to obtain the estimate

$$\begin{aligned}
& \frac{16\kappa_0 B_0^2}{25} \frac{(\log \log n)^2 (\log n)^{\frac{2}{5}}}{n} (1 + \log(K_{\max}(n)q_{\max}(n))) \\
& = O \left( \log n \left( \frac{\log n}{n} \right)^{\frac{2p}{2p+k+1}} \right).
\end{aligned}$$

This implies (5.12). (5.13) is proved similarly observing that as in the proof of the approximation error in (5.9) we obtain

$$\inf_{\alpha \in \mathcal{F}_{n,(K_n,p)}} |\Lambda(\alpha) - \Lambda(\alpha_0)| = O \left( \left( \frac{1}{K_n} \right)^{2p} \right). \quad (5.19)$$

□

The truncation constants  $\beta_n = \frac{1}{5} \log \log n$  are not meant as proposals for practical examples. Note that the same estimates in part b) of the theorem hold for  $\beta_n$  of the form  $cl_n$ , where  $c$  is some bigger constant and  $l_n$  grows slower than  $\log \log n$ .

## 5.2 Net sieves

In this section we apply our results to obtain convergence rates for the conditional log-hazard function for net sieves under the assumption that the

underlying conditional log-hazard function  $\alpha_0$  has an integral representation of the form

$$\alpha_0(t, x) = \int_{\Theta} \Psi(a_{\vartheta}(t, x)) d\nu(\vartheta) \quad (5.20)$$

where  $a_{\vartheta}$ ,  $\vartheta \in \Theta \subset \mathbb{R}^m$ ,  $x \in \mathcal{X} \subset \mathbb{R}^k$  is a set of sieve defining functions,  $\Psi \circ a_{\vartheta}$  is the continuous net and  $\nu$  is a signed measure of bounded variation on  $\Theta$ . This kind of representation is typically related to some smoothness classes (see Yukich, Stichcombe, and White (1995)). Some approximation results by finite nets with rates of approximation are given in Döhler and Rüschendorf (2001a) and applied in the following. Let  $\mathcal{Z}_{k+1}(\mathcal{F}_0)$  denote the class of all functions satisfying (5.20).

Define the basis of the net  $\mathcal{F}_0 = \{\Psi \circ a_{\vartheta}; \vartheta \in \Theta\}$  and for  $\beta > 0$ ,  $K \in \mathbb{N}$  and the finite approximation net

$$\mathcal{F}(\beta, K) = \left\{ \alpha : \mathcal{T} \times \mathcal{X} \rightarrow \mathbb{R}; \alpha(t, x) = \sum_{i=1}^K c_i f_i(t, x), f_i \in \mathcal{F}_0, \sum_{i=1}^K |c_i| \leq \beta \right\}. \quad (5.21)$$

The following conditions were introduced in Döhler and Rüschendorf (2001a) to prove approximation rates by finite nets. Let  $\mu$  be a probability measure on  $\mathbb{R}^{d+1}$ .

**A1)** There exists a  $D > 1$  such that

$$N(\delta, \mathcal{F}_0, d_{L^2(\mu)}) = O\left(\left(\frac{1}{\delta}\right)^{2(D-1)}\right). \quad (5.22)$$

**A2)** Define  $b_z(\vartheta) = a_{\vartheta}(z)$ ,  $z \in \mathcal{T} \times \mathcal{X}$ , then the class  $\{\Psi \circ b_z, z \in \mathbb{R}^{k+1}\}$  is a  $P$ -Donsker class for any probability measure  $P$  on  $\Theta$ .

**Theorem 5.4** *Assume conditions A1) and A2). Let  $K_n := n^{\frac{1}{2+\frac{1}{D-1}}}$ ,  $\beta_n := \frac{1}{5} \log \log n$ , and consider the net ML-estimator*

$$\hat{\alpha}_n := \arg \max_{\alpha \in \mathcal{F}(\beta_n, K_n)} L_n(\alpha).$$

*Then for any  $\alpha_0 \in \mathcal{Z}_{k+1}(\Psi, \mathcal{T}, \mathcal{T} \times \mathcal{X})$  holds*

$$E \|\hat{\alpha}_n - \alpha_0\|_2^2 = O\left(\frac{(\log n)^2}{n^{\frac{1}{2} + \frac{1}{4D-2}}}\right) \quad (5.23)$$

*and*

$$E |\Lambda(\hat{\alpha}_n) - \Lambda(\alpha_0)| = O\left(\frac{(\log n)^2}{n^{\frac{1}{2} + \frac{1}{4D-2}}}\right). \quad (5.24)$$

**Proof:** We apply Theorem 4.1. The approximation error was estimated in Döhler and Rüschendorf (2001a). Let  $\nu_{\alpha_0}$  be the signed measure representing  $\alpha_0$ , then for  $n$  with  $\beta_n \geq 2|\nu_{\alpha_0}|$  holds

$$\inf_{\alpha \in \mathcal{F}(\beta_n, K_n)} \|\alpha - \alpha_0\|_2^2 = O\left(\left(\frac{1}{K_n}\right)^{1+\frac{1}{D-1}}\right). \quad (5.25)$$

Next we prove that for  $\beta > 0$ ,  $K \in \mathbb{N}$  holds:

$$\mathcal{C}_n(\mathcal{F}(\beta, K)) \leq C(D)^K (\beta K)^{2K(2D-1)} n^{2K(2D-1)}. \quad (5.26)$$

Define  $\mathcal{F}'(\beta, K) := \{\alpha : \mathcal{T} \times \mathcal{X} \rightarrow \mathbb{R}, (t, x) \mapsto \sum_{i=1}^K c_i f_i(t, x) \mid f_i \in \mathcal{F}_0, |c_i| \leq \beta\}$ . Then  $\mathcal{F}(\beta, K) \subset \mathcal{F}'(\beta, K)$ , and we obtain for any probability measure  $\nu$  on  $\mathcal{T} \times \mathcal{X}$  and  $\delta > 0$  using some well-known rules for covering numbers (cf. van der Vaart and Wellner (1996)).

$$\begin{aligned} & N(\delta, \mathcal{F}'(\beta, K), d_{L^1(\nu)}) \\ & \leq N(\delta, \mathcal{F}'(\beta, K), d_{L^2(\nu)}) \\ & \leq N\left(\delta, 10s \bigoplus_{i=1}^K [-\beta, \beta] \odot \mathcal{F}_0, d_{L^2(\nu)}\right) \leq N\left(\frac{\delta}{K}, [-\beta, \beta] \odot \mathcal{F}_0, d_{L^2(\nu)}\right)^K \\ & \leq \left[N\left(\frac{\delta}{2\beta K}, \mathcal{F}_0, d_{L^2(\nu)}\right) \frac{4\beta K}{\delta}\right]^K \leq \left[C\left(\frac{\delta}{2\beta K}\right)^{-2(D-1)} \frac{4\beta K}{\delta}\right]^K \\ & = C(D)^K (\beta K)^{K(2D-1)} \left(\frac{1}{\delta}\right)^{K(2D-1)} \end{aligned}$$

independent of  $\nu$  as in (5.26).

From Theorem 4.1 with  $M = \beta_n$  this implies

$$\begin{aligned} & E\|\hat{\alpha}_n - \alpha_0\|_2^2 \\ & = O\left(\exp(3\beta_n) \left[\left(\frac{1}{K_n}\right)^{1+\frac{1}{D-1}} + \frac{K_n}{n} \log(C(D)(n\beta_n K_n)^{2(2D-1)})\right]\right) \\ & = O\left(\log n \left[\left(\frac{1}{n}\right)^{\frac{1+\frac{1}{D-1}}{2+\frac{1}{D-1}}} + \left(\frac{1}{n}\right)^{1-\frac{1}{2+\frac{1}{D-1}}} \log n\right]\right) \\ & = O\left((\log n)^2 \left(\frac{1}{n}\right)^{1-\frac{1}{2+\frac{1}{D-1}}}\right) \\ & = O\left((\log n)^2 \left(\frac{1}{n}\right)^{\frac{1}{2}+\frac{1}{4D-2}}\right), \end{aligned}$$

and the result follows.

The proof of (5.24) is analogous using the approximation estimate

$$\inf_{\alpha \in \mathcal{F}(\beta_n, K_n)} |\Lambda(\alpha) - \Lambda(\alpha_0)| = O \left( \left( \frac{1}{K_n} \right)^{1 + \frac{1}{D-1}} \right). \quad (5.27)$$

For the proof note that for  $n$  with  $\beta_n \geq L := \max\{2|\nu_{\alpha_0}|, \|\alpha_0\|_\infty\}$  it holds by (1.8) that

$$\begin{aligned} \inf_{\alpha \in \mathcal{F}(\beta_n, K_n)} |\Lambda(\alpha) - \Lambda(\alpha_0)| &\leq \inf_{\alpha \in \mathcal{F}(L, K_n)} |\Lambda(\alpha) - \Lambda(\alpha_0)| \\ &\leq k'(L) \inf_{\alpha \in \mathcal{F}(L, K_n)} \|\alpha - \alpha_0\|_2^2 \\ &= O \left( \left( \frac{1}{K_n} \right)^{1 + \frac{1}{D-1}} \right); \end{aligned}$$

the last estimate is from Döhler and Rüschendorf (2001a).  $\square$

As in section 5.1 alternative choices of the truncation constants  $\beta_n$  are possible.

We consider the special classes of neural nets, wavelet nets, and radial basis function nets. In the following examples we use some approximation results from Döhler and Rüschendorf (2001a).

- a) **Neural nets** Here  $\mathcal{F}_0 = \{f_0 : \mathcal{T} \times \mathcal{X} \rightarrow [0, 1], z \mapsto \Psi(\gamma z + \delta) \mid \gamma \in \mathbb{R}^{k+1}, \delta \in \mathbb{R}\}$ , where  $\Psi : \mathbb{R} \rightarrow [0, 1]$  is of bounded variation. Then the conditions A1) and A2) are fulfilled with  $D = k + 4$  and Theorem 5.4 implies

$$E \|\hat{\alpha}_n - \alpha_0\|_2^2 = O \left( \frac{(\log n)^2}{n^{\frac{1}{2} + \frac{1}{4k+14}}} \right). \quad (5.28)$$

The same rate holds if the representation property of  $\alpha_0$  is replaced by Barron's (1993) finiteness condition on the Fourier transform

$$C_f = \int |w|_1 |\tilde{f}(w)| dw < \infty \quad (5.29)$$

where  $\tilde{f}$  is the Fourier transform of  $f$ .

If  $P^{(T, X)}$  has a density with bounded support the convergence rate can be improved to

$$E \|\hat{\alpha}_n - \alpha_0\|_2^2 = O \left( \frac{(\log n)^2}{n^{\frac{1}{2} + \frac{1}{4k+6}}} \right). \quad (5.30)$$

Similar rates with  $\frac{1}{2}$  instead of  $\frac{1}{2} + \frac{1}{4k+6}$  in the exponent were obtained previously for regression estimation, in Barron (1994) and for density estimation in Modha and Masry (1996).

- b) **Radial basis-function nets** Here  $\mathcal{F}_0 = \{f_0 : \mathcal{T} \times \mathcal{X} \rightarrow [0, 1], z \mapsto \varrho(\|\gamma(z - \delta)\|) \mid \gamma \in \mathbb{R}^{k+1}, \delta \in \mathbb{R}\}$ , where  $\varrho : \mathbb{R}^+ \rightarrow [0, 1]$  is monotonically non-increasing. Then the conditions A1) and A2) are fulfilled with  $D = k + 5$  and from Theorem 5.4 we obtain

$$E\|\hat{\alpha}_n - \alpha_0\|_2^2 = O\left(\frac{(\log n)^2}{n^{\frac{1}{2} + \frac{1}{4k+18}}}\right). \quad (5.31)$$

- c) **Wavelet nets** Here  $\mathcal{F}_0 = \{f_0 : \mathcal{T} \times \mathcal{X} \rightarrow [0, 1], z \mapsto \Psi(\gamma(z - \delta)) \mid \gamma \in \mathbb{R}^{k+1}, \delta \in \mathbb{R}\}$ , where  $\Psi : \mathbb{R}^{d+1} \rightarrow [0, 1]$  is Lipschitz with bounded support. Then by Theorem 5.4 with  $D = 3k + 4$  we obtain

$$E\|\hat{\alpha}_n - \alpha_0\|_2^2 = O\left(\frac{(\log n)^2}{n^{\frac{1}{2} + \frac{1}{12k+14}}}\right). \quad (5.32)$$

Note that in all three cases a corresponding convergence result also holds in terms of the  $\Lambda$ -distance.

## Résumé

In conclusion this paper gives quite general results on the convergence rates for sieved minimum contrast estimators and also for the related adaptive versions of these estimators. The results are formulated in detail for the example of estimating the log-hazard function in censoring models. In comparison to the related general approach in Birgé and Massart (1998) and Barron, Birgé, and Massart (1999) we use some simpler conditions concerning the covering numbers. The results in this paper are illustrated with examples of sieves such as neural nets, wavelet nets, radial basis function nets and tensor product splines. Some further applications of the method in this paper to more general type of censorings as well as to a more detailed study of neural net estimators is given in the forthcoming papers in Döhler and Rüschendorf (2000, 2001b).

## References

- Andersen, P. K., O. Borgan, R. Gill, and N. Keiding (1993). *Statistical models based on counting processes*. Springer.
- Barron, A. R. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory* 39(3), 930–945.
- Barron, A. R. (1994). Approximation and estimation bounds for artificial neural networks. *Machine Learning* 14, 115–133.

- Barron, A. R., L. Birgé, and P. Massart (1999). Risk bounds for model selection via penalization. *Probability Theory and Related Fields* 113(3), 301–413.
- Birgé, L. and P. Massart (1998). Minimum contrast estimators on sieves: Exponential bounds and rates of convergence. *Bernoulli* 4(3), 329–375.
- Döhler, S. (2000a). Consistent hazard regression estimation by sieved maximum likelihood estimators. In *Proceedings of Conference on Limit Theorems in Balatonlelle (to appear)*.
- Döhler, S. (2000b). *Empirische Risiko-Minimierung bei zensierten Daten*. Dissertation, Universität Freiburg, <http://webdoc.sub.gwdg.de/ebook/e/2001/freidok/69.pdf>.
- Döhler, S. and L. Rüschendorf (2000). A consistency result in general censoring models. *Statistics (to appear)*.
- Döhler, S. and L. Rüschendorf (2001a). An approximation result for nets in functional estimation. *Statistics & Probability Letters* 52, 373–380.
- Döhler, S. and L. Rüschendorf (2001b). On adaptive estimation by neural net type estimators. In *Proceedings of MSRI Workshop on Nonlinear Estimation and Classification (to appear)*, Berkeley 2001.
- Kohler, M. (1997). *Nichtparametrische Regressionsschätzung mit Splines*. Dissertation, Universität Stuttgart, <http://www.mathematik.uni-stuttgart.de/mathA/lst3/kohler/papers!html>.
- Kohler, M. (1999a). Nonparametric estimation of piecewise smooth regression functions. *Statistics & Probability Letters* 43, 49–55.
- Kohler, M. (1999b). Universally consistent regression function estimation using hierarchical B-splines. *Journal of Multivariate Analysis* 68, 138–164.
- Kooperberg, C., C. J. Stone, and Y. K. Truong (1995a). Hazard regression. *Journal of the American Statistical Association* 90, 78–94.
- Kooperberg, C., C. J. Stone, and Y. K. Truong (1995b). The  $L_2$  rate of convergence for hazard regression. *Scandinavian Journal of Statistics* 22, 143–157.
- Krzyzak, A. and T. Linder (1998). Radial basis function networks and computational regularization in function learning. *IEEE Transactions on Information Theory* 9, 247–256.
- Lee, W., P. Bartlett, and R. Williamson (1996). Efficient agnostic learning of neural networks with bounded fan-in. *IEEE Transactions on Information Theory* 42(6), 2118–2132.
- Lugosi, G. and K. Zeger (1995). Nonparametric estimation via empirical risk minimization. *IEEE Transactions on Information Theory* 41(3), 677–687.



- Modha, D. and E. Masry (1996). Rate of convergence in density estimation using neural networks. *Neural Computation* 8, 1107–1122.
- Pollard, D. (1984). *Convergence of stochastic processes*, Volume XIV of *Series in Statistics*. Springer.
- Pollard, D. (1990). *Empirical Processes: Theory and Applications*. Institute of Mathematical Statistics, Hayward.
- van der Vaart, A. and J. Wellner (1996). *Weak convergence and empirical processes*. Springer, New York.
- van Keilegom, I. and N. Veraverbeke (2001). Hazard rate estimation in non-parametric regression with censored data. *Annals of the Institute of Statistical Mathematics* 53, 730–745.
- Vapnik, V. (1995). *The Nature of Statistical Learning Theory*. Springer, New York.
- Wong, W. and X. Shen (1995). Probability inequalities for likelihood ratios and convergence rates of sieve mles. *Annals of Statistics* 23, 339–362.
- Yang, Y. and A. Barron (1998). An asymptotic property of model selection criteria. *IEEE Transactions on Information Theory* 44, 95–116.
- Yukich, J., M. Stichcombe, and H. White (1995). Sup-norm approximation bounds for networks through probabilistic methods. *IEEE Transactions on Information Theory* 41(4), 1021–1027.

corresponding author:

Ludger Rüschendorf  
 University of Freiburg  
 Institute for Mathematical Stochastics  
 Eckerstr. 1  
 79104 Freiburg  
 Germany

Tel.: +49-761-2035665  
 Fax: +49-761-2035661

E-mail: [ruschen@stochastik.uni-freiburg.de](mailto:ruschen@stochastik.uni-freiburg.de)