



Published in final edited form as:

*J Multivar Anal.* 2009 March 1; 100(3): 345–362.

## The Penalized Profile Sampler

GUANG CHENG<sup>a,\*</sup> and MICHAEL R. KOSOROK<sup>b</sup>

<sup>a</sup> Department of Statistical Science, Duke University, 214 Old Chemistry Building, Durham, NC 27708, USA

<sup>b</sup> Department of Biostatistics, School of Public Health, University of North Carolina at Chapel Hill, 3101 McGavran-Greenberg Hall, Chapel Hill, NC 27599, USA

### Abstract

The penalized profile sampler for semiparametric inference is an extension of the profile sampler method [9] obtained by profiling a penalized log-likelihood. The idea is to base inference on the posterior distribution obtained by multiplying a profiled penalized log-likelihood by a prior for the parametric component, where the profiling and penalization are applied to the nuisance parameter. Because the prior is not applied to the full likelihood, the method is not strictly Bayesian. A benefit of this approximately Bayesian method is that it circumvents the need to put a prior on the possibly infinite-dimensional nuisance components of the model. We investigate the first and second order frequentist performance of the penalized profile sampler, and demonstrate that the accuracy of the procedure can be adjusted by the size of the assigned smoothing parameter. The theoretical validity of the procedure is illustrated for two examples: a partly linear model with normal error for current status data and a semiparametric logistic regression model. Simulation studies are used to verify the theoretical results.

### Keywords

Penalized Likelihood; Posterior Distribution; Profile Likelihood; Semiparametric Inference; Smoothing Parameter

## 1 Introduction

Semiparametric models are statistical models indexed by both a finite dimensional parameter of interest  $\theta$  and an infinite dimensional nuisance parameter  $\eta$ . In order to make statistical inference about  $\theta$  separately from  $\eta$ , we estimate the nuisance parameter with  $\hat{\eta}_\theta$ , its maximum likelihood estimate at each fixed  $\theta$ , i.e.

$$\hat{\eta}_\theta = \operatorname{argmax}_{\eta \in \mathcal{H}} \operatorname{lik}_n(\theta, \eta),$$

where  $\operatorname{lik}_n(\theta, \eta)$  is the likelihood of the semiparametric model given  $n$  observations and  $\mathcal{H}$  is the parameter space for  $\eta$ . Therefore we can do frequentist inference about  $\theta$  based on the profile likelihood, which is typically defined as

\*Corresponding author. Department of Statistical Science, DUKE University, 214 Old Chemistry Building, Durham, NC 27708, USA. Fax: 919 684 8594. chengg@duke.edu (G. Cheng), kosorok@unc.edu (M.R. Kosorok).

$$p l_n(\theta) = \sup_{\eta \in \mathcal{H}} \text{lik}_n(\theta, \eta).$$

The convergence rate of the nuisance parameter  $\eta$  is the order of  $d(\hat{\eta}_{\hat{\theta}_n}, \eta_0)$ , where  $d(\cdot, \cdot)$  is some metric on  $\eta$ ,  $\hat{\theta}_n$  is any sequence satisfying  $\hat{\theta}_n = \theta_0 + o_P(1)$ , and  $(\eta_0, \theta_0)$  is the true value of  $(\eta, \theta)$ . Typically,

$$d(\hat{\eta}_{\hat{\theta}_n}, \eta_0) = O_p(\|\hat{\theta}_n - \theta_0\| + n^{-r}), \quad (1)$$

where  $\|\cdot\|$  is the Euclidean norm and  $r > 1/4$ . Of course, a smaller value of  $r$  leads to a slower convergence rate of the nuisance parameter. For instance, the nuisance parameter in the Cox proportional hazards model with right censored data, the cumulative hazard function, has the parametric rate, i.e.,  $r = 1/2$ . If current status data is applied to the Cox model instead, then the convergence rate will be slower, with  $r = 1/3$ , due to the loss of information provided by this kind of data.

The profile sampler is the procedure of sampling from the posterior of the profile likelihood in order to estimate and draw inference on the parametric component  $\theta$  in a semiparametric model, where the profiling is done over the possibly infinite-dimensional nuisance parameter  $\eta$ . [9] show that the profile sampler gives a first order correct approximation to the maximum likelihood estimator  $\hat{\theta}_n$  and consistent estimation of the efficient Fisher information for  $\theta$  even when the nuisance parameter is not estimable at the  $\sqrt{n}$  rate. Another Bayesian procedure employed to do semiparametric estimation is considered in [17] who study the marginal semiparametric posterior distribution for a parameter of interest. In particular, [17] show that marginal semiparametric posterior distributions are asymptotically normal and centered at the corresponding maximum likelihood estimates or posterior means, with covariance matrix equal to the inverse of the Fisher information. Unfortunately, this fully Bayesian method requires specification of a prior on  $\eta$ , which is quite challenging since for some models there is no direct extension of the concept of a Lebesgue dominating measure for the infinite-dimensional parameter set involved [8]. The advantages of the profile sampler for estimating  $\theta$  compared to other methods is discussed extensively in [2], [3] and [9].

The motivation for studying second order asymptotic properties of the profile sampler comes from the observed simulation differences in the Cox model with different types of data, i.e. right censored data [2] and current status data [9]. The profile sampler generated based on the first model yields much more accurate estimation results comparing to the second model when the sample size is relatively small. [2] and [3] have successfully explored the theoretical reasons behind the above phenomena by establishing the relation between the estimation accuracy of the profile sampler, measured in terms of second order asymptotics, and the convergence rate of the nuisance parameters. Specifically speaking, the profile sampler generated from a semiparametric model with a faster convergence rate usually yields more precise frequentist inference of  $\theta$ . These second order results are verified in [2] and [3] for several examples, including the proportional odds model, case-control studies with missing covariates, and the partly linear model. The convergence rates for these models range from the parametric to the cubic. The work in [3] has shown clearly that the accuracy of the inference for  $\theta$  based on the profile sampler method is intrinsically determined by the semiparametric model specifications through its entropy number.

In many semiparametric models involving a smooth nuisance parameter, it is often convenient and beneficial to perform estimation using penalization. One motivation for this is that, in the absence of any restrictions on the form of the function  $\eta$ , maximum likelihood estimation for some semiparametric models leads to over-fitting. Seminal applications of penalized maximum likelihood estimation include estimation of a probability density function in [18] and nonparametric linear regression in [19]. Note that penalized likelihood is a special case of penalized quasi-likelihood studied in [13]. Under certain reasonable regularity conditions, penalized semiparametric log-likelihood estimation can yield fully efficient estimates for  $\theta$  (see, for example, [13]). As far as we are aware, the only general procedure for inference for  $\theta$  in this context known to be theoretically valid is a weighted bootstrap with bounded random weights (see [11]). It is even unclear whether the usual nonparametric bootstrap will work in this context when the nuisance parameter has a convergence rate  $r < 1/2$ .

The purpose of this paper is to ask the somewhat natural question: does sampling from the exponential of a profiled penalized log-likelihood (which process we refer hereafter to as the penalized profile sampler) yield first and even second order accurate frequentist inference? The conclusion of this paper is that the answer is yes and, moreover, the accuracy of the inference depends in a fairly simple way on the size of the smoothing parameter.

The unknown parameters in the semiparametric models we study in this paper include  $\theta$ , which we assume belongs to some compact set  $\Theta \subset \mathbb{R}^d$ , and  $\eta$ , which we assume to be a function in the Sobolev class of functions  $\mathcal{H}_k$  or its subset  $\mathcal{H}_k^M \equiv \mathcal{H}_k \cap \{\eta: \|\eta\|_\infty \leq M\}$  for some known  $M < \infty$  supported on some compact set on the real line. The Sobolev class of functions  $\mathcal{H}_k$  is defined as the set  $\{\eta: J^2(\eta) \equiv \int_{\mathcal{Z}} (\eta^{(k)}(z))^2 dz < \infty\}$ , where  $\eta^{(j)}$  is the  $j$ -th derivative of  $\eta$  with respect to  $z$ . Obviously  $J^2(\eta)$  is some measurement of complexity of  $\eta$ . We denote  $\mathcal{H}_k$  as the Sobolev function class with degree  $k$ . The penalized log-likelihood in this context is:

$$\log lik_{\lambda_n}(\theta, \eta) = \log lik(\theta, \eta) - n\lambda_n^2 J^2(\eta), \tag{2}$$

where  $\log lik(\theta, \eta) \equiv n\mathbb{P}_n \ell_{\theta, \eta}(X)$ ,  $\ell_{\theta, \eta}(X)$  is the log-likelihood of the single observation  $X$ , and  $\lambda_n$  is a smoothing parameter, possibly dependent on data. In practice,  $\lambda_n$  can be obtained by cross-validation [23] or by inspecting the various curves for different values of  $\lambda_n$ . The penalized maximum likelihood estimators  $\hat{\theta}_n$  and  $\hat{\eta}_n$  depend on the choice of the smoothing parameter  $\lambda_n$ . Consequently we use the notation  $\hat{\theta}_{\lambda_n}$  and  $\hat{\eta}_{\lambda_n}$  for the remainder of this paper to denote the estimators obtained from maximizing (2). In particular, a larger smoothing parameter usually leads to a less rough penalized estimator of  $\eta_0$ . It is of interest to establish the asymptotic property of the proposed penalized profile sampler procedure with a data-driven  $\lambda_n$ . Further studies on this issue are needed, but it is beyond the scope of this paper.

For the purpose of establishing first order accuracy of inference for  $\theta$  based on the penalized profile sampler, we assume that the bounds for the smoothing parameter are in the form below:

$$\lambda_n = o_p(n^{-1/4}) \text{ and } \lambda_n^{-1} = O_p(n^{k/(2k+1)}). \tag{3}$$

The condition (3) is assumed to hold throughout this paper. One way to ensure (3) in practice is simply to set  $\lambda_n = n^{-k/(2k+1)}$ . Or we can just choose  $\lambda_n = n^{-1/3}$  which is independent of  $k$ . It turns out that the upper bound guarantees that  $\hat{\theta}_{\lambda_n}$  is  $\sqrt{n}$ -consistent, while the lower bound controls the penalized nuisance parameter estimator convergence rate. Another approach to controlling estimators is to use sieve estimates with assumptions on the derivatives (see [6]). We will not pursue this further here.

The log-profile penalized likelihood is defined as follows:

$$\log pl_{\lambda_n}(\theta) = \log lik(\theta, \hat{\eta}_{\theta, \lambda_n}) - n\lambda_n^2 J^2(\hat{\eta}_{\theta, \lambda_n}), \tag{4}$$

where  $\hat{\eta}_{\theta, \lambda_n}$  is  $\operatorname{argmax}_{\eta \in \mathcal{A}_k} \log lik_{\lambda_n}(\theta, \eta)$  for fixed  $\theta$  and  $\lambda_n$ . Note that  $J(\hat{\eta}_{\tilde{\theta}_n, 0}) \geq J(\hat{\eta}_{\tilde{\theta}_n, \lambda_n})$ , where  $\eta_{\theta, 0} = \hat{\eta}_{\theta} \equiv \operatorname{argmax}_{\eta \in \mathcal{A}_k} \log lik(\theta, \eta)$  for a fixed  $\theta$ , based on the inequality that  $\log lik_{\lambda_n}(\theta, \hat{\eta}_{\tilde{\theta}_n, 0}) \leq \log lik_{\lambda_n}(\theta, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})$ . Hence again we verify that the smoothing parameter  $\lambda_n$  plays a role in determining the complexity degree of the estimated nuisance parameter. The penalized profile sampler is just the procedure of sampling from the posterior distribution of  $pl_{\lambda_n}(\theta)$  by assigning a prior on  $\theta$ . By analyzing the corresponding MCMC chain from the frequentist's point of view, our paper obtains the following conclusions:

1. *Distribution Approximation:* The posterior distribution with respect to  $pl_{\lambda_n}(\theta)$  can be approximated by the normal distribution with mean the maximum penalized likelihood estimator of  $\theta$  and variance the inverse of the efficient information matrix, with error  $O_p(n^{1/2}\lambda_n^2)$ ;
2. *Moment Approximation:* The maximum penalized likelihood estimator of  $\theta$  can be approximated by the mean of the MCMC chain with error  $O_p(\lambda_n^2)$ . The efficient information matrix can be approximated by the inverse of the variance of the MCMC chain with error  $O_p(n^{1/2}\lambda_n^2)$ ;
3. *Confidence Interval Approximation:* An exact frequentist confidence interval of Wald's type for  $\theta$  can be estimated by the credible set obtained from the MCMC chain with error  $O_p(\lambda_n^2)$ .

Obviously, given any smoothing parameter satisfying the upper bound in (3), the penalized profile sampler can yield first order frequentist valid inference for  $\theta$ , similar as to what was shown for the profile sampler in [9]. Moreover, the above conclusions are actually second order frequentist valid results, whose approximation accuracy is directly controlled by the smoothing parameter. Note that the corresponding results for the usual (non-penalized) profile sampler with nuisance parameter convergence rate  $r$  in [3] are obtained by replacing in the above  $O_p(n^{1/2}\lambda_n^2)$  with  $O_p(n^{-1/2} \vee n^{-r+1/2})$  and  $O_p(\lambda_n^2)$  with  $O_p(n^{-1} \vee n^{-r})$ , for all respective occurrences, where  $r$  is as defined in (1).

Our results are the first general higher order frequentist inference results for penalized semi-parametric estimation. We also note, however, that some results on second order efficiency of semiparametric estimators were derived in [4]. The layout of the article is as follows. The next section, section 2, introduces the two main examples we will be using for illustration: partly linear regression for current status data and semiparametric logistic regression. Some background is given in section 3, including the concept of a least favorable submodel as well as the main model assumptions. One preliminary theorem concerning about second order asymptotic expansions of the log-profile penalized likelihood is also presented in section 3. The main results and implications are discussed in section 4, and all remaining model assumptions are verified for the examples in section 5. A brief discussion of future work is given in section 6. We postpone all technical tools and proofs to the last section, section 7.

## 2 Examples

### 2.1 Partly Linear Normal Model with Current Status Data

In this example, we study the partly linear regression model with normal residue error. The continuous outcome  $Y$ , conditional on the covariates  $(U, V) \in \mathbb{R}^d \times \mathbb{R}$ , is modeled as

$$Y = \theta^T U + f(V) + \varepsilon, \tag{5}$$

where  $f$  is an unknown smooth function, and  $\varepsilon \sim N(0, \sigma^2)$  with finite variance  $\sigma^2$ . For simplicity, we assume for the rest of the paper that  $\sigma = 1$ . The theory we propose also works when  $\sigma$  is unknown, but the added complexity would detract from the main issues. We also assume that only the current status of response  $Y$  is observed at a random censoring time  $C \in \mathbb{R}$ . In other words, we observe  $X = (C, \Delta, U, V)$ , where indicator  $\Delta = 1\{Y \leq C\}$ . Current status data may occur due to study design or measurement limitations. Examples of such data arise in several fields, including demography, epidemiology and econometrics. For simplicity of exposition,  $\theta$  is assumed to be one dimensional.

Under the model (5) and given that the joint distribution for  $(C, U, V)$  does not involve parameters  $(\theta, f)$ , the log-likelihood for a single observation at  $X = x \equiv (c, \delta, u, v)$  is

$$\text{loglik}_{\theta, f}(x) = \delta \log\{\Phi(c - \theta u - f(v))\} + (1 - \delta) \log\{1 - \Phi(c - \theta u - f(v))\}, \tag{6}$$

where  $\Phi$  is the cdf of the standard normal distribution. The parameter of interest,  $\theta$ , is assumed to belong to some compact set in  $\mathbb{R}^1$ . The nuisance parameter is the function  $f$ , which belongs to the Sobolev function class of degree  $k$ . We further make the following assumptions on this model. We assume that  $(Y, C)$  is independent given  $(U, V)$ . The covariates  $(U, V)$  are assumed to belong to some compact set, and the support for random censoring time  $C$  is an interval  $[l_c, u_c]$ , where  $-\infty < l_c < u_c < \infty$ . In addition,  $P\text{Var}(U|V)$  is strictly positive and  $Pf(V) = 0$ . The first order asymptotic behaviors of the penalized log-likelihood estimates of a slightly more general version of this model have been extensively studied in [10].

### 2.2 Semiparametric Logistic Regression

Let  $X_1 = (Y_1, W_1, Z_1), X_2 = (Y_2, W_2, Z_2), \dots$  be independent copies of  $X = (Y, W, Z)$ , where  $Y$  is a dichotomous variable with conditional expectation  $P(Y|W, Z) = F(\theta^T W + \eta(Z))$ .  $F(u)$  is the logistic distribution defined as  $e^u / (e^u + 1)$ . Obviously the likelihood for a single observation is of the following form:

$$\text{lik}_{\theta, \eta}(x) = F(\theta^T w + \eta(z))^y (1 - F(\theta^T w + \eta(z)))^{1-y} f^{(W,Z)}(w, z). \tag{7}$$

This example is a special case of quasi-likelihood in partly linear models when the conditional variance of response  $Y$  is taken to have some quadratic form of the conditional mean of  $Y$ . In the absence of any restrictions on the form of the function  $\eta$ , the maximum likelihood of this simple model often leads to over-fitting. Hence [5] propose maximizing instead the penalized likelihood of the form  $\text{loglik}(\theta, \eta) - n\lambda_n^2 J^2(\eta)$ ; and [13] showed the asymptotic consistency of the maximum penalized likelihood estimators for  $\theta$  and  $\eta$ . For simplicity, we will restrict ourselves to the case where  $\Theta \subset \mathbb{R}^1$  and  $(W, Z)$  have bounded support, say  $[0, 1]^2$ . To ensure the identifiability of the parameters, we assume that  $P\text{Var}(W|Z)$  is positive and that the support of  $Z$  contains at least  $k$  distinct points in  $[0, 1]$ , see lemma 7.1 in [15].

**Remark 1**—Another interesting potential example we may apply the penalized profile sampler method to is the classic proportional hazards model with current status data by penalizing the cumulative hazard function with its Sobolev norm. There are two motivations for us to penalize the cumulative hazard function in the Cox model. One is that the estimated

step functions from the unpenalized estimation cannot be used easily for other estimation or inference purposes. Another issue with the unpenalized approach is that without making stronger continuity assumptions, we cannot achieve uniform consistency even on a compact set [10]. The asymptotic properties of the corresponding penalized M-estimators have been studied in [12].

### 3 Preliminaries

In this section, we present some necessary preliminary material concerning least favorable submodels and assume some structural requirements to achieve second order asymptotic expansion of the log-profile penalized likelihood (21).

#### 3.1 Least favorable submodels

In this subsection, we briefly review the concept of a least favorable submodel. A submodel  $t \mapsto \text{lik}_{t,\eta_t}$  is defined to be least favorable at  $(\theta, \eta)$  if  $\tilde{\ell}_{\theta,\eta} = \partial/\partial t \log \text{lik}_{t,\eta_t}$ , given  $t = \theta$ , where  $\tilde{\ell}_{\theta,\eta}$  is the efficient score function for  $\theta$ . The efficient score function for  $\theta$  can be viewed as the projection of the score function for  $\theta$  onto the tangent space of  $\eta$ . The inverse of its variance is exactly the efficient information matrix  $\tilde{I}_{\theta,\eta}$ . We abbreviate hereafter  $\tilde{\ell}_{\theta_0\eta_0}$  and  $\tilde{I}_{\theta_0,\eta_0}$  with  $\tilde{\ell}_0$  and  $\tilde{I}_0$ , respectively. The “direction” along which  $\eta_t$  approaches  $\eta$  in the least favorable submodel is called the least favorable direction. An insightful review about least favorable submodels and efficient score functions can be found in Chapter 3 of [7]. We assume that in our setting a least favorable submodel always exists. By the above construction of the least favorable submodel,  $\log pl_{\lambda_n}(\theta)$  can be rewritten in the following form:

$$\log pl_{\lambda_n}(\theta) = n(\mathbb{P}_n \ell(\theta, \theta, \tilde{\eta}_{\theta, \lambda_n}) - \lambda_n^2 J^2(\eta_\theta(\theta, \tilde{\eta}_{\theta, \lambda_n}))), \tag{8}$$

where  $\ell(t, \theta, \eta)(x) = \ell_{t,\eta_t(\theta,\eta)}(x)$ ,  $t \mapsto \eta_t(\theta, \eta)$  is a general map from the neighborhood of  $\theta$  into the parameter set for  $\eta$ , with  $\eta_\theta(\theta, \eta) = \eta$ . The concrete forms of (8) will depend on the situation.

The derivatives of the function  $\ell(t, \theta, \eta)$  are with respect to its first argument,  $t$ . For the derivatives relative to the argument  $\theta$ , we use the following shortened notation:  $\ell_\theta(t, \theta, \eta)$  indicates the first derivative of  $\ell(t, \theta, \eta)$  with respect to  $\theta$  and  $\ell_{t,\theta}(t, \theta, \eta)$  denotes the derivative of  $\ell(t, \theta, \eta)$  relative to  $\theta$ . Also,  $\ell_{t,t}(\theta)$  and  $\ell_{t,\theta}(\eta)$  indicate the maps  $\theta \mapsto \tilde{\ell}(t, \theta, \eta)$  and  $\eta \mapsto \tilde{\ell}_{t,\theta}(t, \theta, \eta)$ , respectively. For brevity, we denote  $\ell_0 = \ell(\theta_0, \theta_0, \eta_0)$ ,  $\tilde{\ell}_0 = \tilde{\ell}(\theta_0, \theta_0, \eta_0)$  and  $\ell_0^{(3)} = \ell^{(3)}(\theta_0, \theta_0, \eta_0)$ , where  $\theta_0$  and  $\eta_0$  are the true values of  $\theta$  and  $\eta$ . Of course, we can write  $\ell_0(X)$  as  $\tilde{\ell}_0(X)$  based on the construction of the least favorable submodel. All the necessary derivatives of  $\ell(t, \theta, \eta)$  w.r.t.  $t$  or  $\theta$  in this paper are assumed to have integrable envelope functions in some neighborhood of  $(\theta_0, \theta_0, \eta_0)$ . In the following, we use  $P_{\theta,\eta}U$  to denote the expectation of a random variable  $U$  at the parameter  $(\theta, \eta)$ , and use  $PU$  to represent  $P_{\theta_0,\eta_0}U$  for simplicity.

#### 3.2 Main Assumptions

The set of structural conditions about the least favorable submodel are the “no-bias” conditions:

$$P \dot{\ell}(\theta_0, \tilde{\theta}_n, \tilde{\eta}_{\tilde{\theta}_n, \lambda_n}) = O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)^2, \tag{9}$$

$$P \ddot{\ell}(\theta_0, \tilde{\theta}_n, \tilde{\eta}_{\tilde{\theta}_n, \lambda_n}) = P \ddot{\ell}_0 + O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|), \tag{10}$$



for any sequence  $\tilde{\theta}_n$  satisfying  $\tilde{\theta}_n = \theta_0 + o_p(1)$ . The verifications of (9) and (10) depend on the smoothness of  $\ell(t, \theta, \eta)$  and the convergence rate of the penalized nuisance parameter based on the functional Taylor expansions around the true values. The convergence rate typically has the following upper bound:

$$d(\widehat{\eta}_{\tilde{\theta}_n, \lambda_n}, \eta_0) = O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|). \quad (11)$$

The form of  $d(\eta, \eta_0)$  may vary for different situations and does not need to be specified in this subsection beyond the given conditions. (11) implies that  $\widehat{\eta}_{\tilde{\theta}_n, \lambda_n}$  is consistent for  $\eta_0$  as  $\tilde{\theta}_n \rightarrow \theta_0$  in probability. Hence (9) and (10) hold provided the Fréchet derivatives of the maps  $\eta \mapsto \ell(\theta_0, \theta_0, \eta)$  and  $\eta \mapsto \ell_{t, \theta}(\theta_0, \theta_0, \eta)$  are bounded, and

$$P \dot{\ell}(\theta_0, \theta_0, \eta) = O(d^2(\eta, \eta_0)), \quad (12)$$

which is usually implied by a bounded Fréchet derivative of  $\eta \mapsto \ell(\theta_0, \theta_0, \eta)$  and second order Fréchet differentiability of the map  $\eta \mapsto \text{lik}(\theta_0, \eta)$ .

The empirical version of the no-bias conditions,

$$\mathbb{P}_n \dot{\ell}(\theta_0, \tilde{\theta}_n, \widehat{\eta}_{\tilde{\theta}_n, \lambda_n}) = \mathbb{P}_n \tilde{\ell}_0 + O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)^2, \quad (13)$$

$$\mathbb{P}_n \ddot{\ell}(\theta_0, \tilde{\theta}_n, \widehat{\eta}_{\tilde{\theta}_n, \lambda_n}) = \mathbb{P} \ddot{\ell}_0 + O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|), \quad (14)$$

where  $\mathbb{P}_n$  represents the empirical distribution of the observations, ensures that the penalized profile likelihood behaves like a penalized likelihood in the parametric model asymptotically and therefore yields a second order asymptotic expansion of the penalized profile log-likelihood. Obviously the empirical no-bias conditions are built upon (9) and (10) by assuming the sizes of the collections of the functions  $\dot{\ell}$  and  $\ddot{\ell}$  are manageable. This condition is expressed in the language of empirical processes. Provided that  $\tilde{\ell}_0$  and  $\ell_{t, \theta}(\theta_0, \theta_0, \eta_0)$  are square integrable, (14) follows from (10) if we assume

$$\mathbb{G}_n(\ddot{\ell}(\theta_0, \tilde{\theta}_n, \widehat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \ddot{\ell}_0) = o_p(1), \quad (15)$$

where  $\mathbb{G}_n \equiv \sqrt{n}(\mathbb{P}_n - P)$  is used for the empirical processes of the observations. If we further assume that

$$\mathbb{G}_n(\ell_{t, \theta}(\theta_0, \tilde{\theta}_n, \widehat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \ell_{t, \theta}(\theta_0, \theta_0, \eta_0)) = o_p(1), \quad (16)$$

$$\mathbb{E}_n(\dot{\ell}(\theta_0, \theta_0, \widehat{\eta}_{\theta_n, \lambda_n}) - \dot{\ell}_0) = O_p(n^{-\frac{1}{4k+2}}(\lambda_n + \|\widetilde{\theta}_n - \theta_0\|)), \tag{17}$$

for any sequence  $\widetilde{\theta}_n$  satisfying  $\widetilde{\theta}_n = \theta_0 + o_p(1)$ , then (13) follows. Note that the conditions (15)–(17) are concerned with the asymptotic equicontinuity of the empirical process measure of  $\dot{\ell}$ ,  $\dot{\ell}_{t,\theta}$  and  $\dot{\ell}$ , respectively. Thus we will be able to use technical tools T2 and T5 given in the appendix to show (15)–(17). We next present the preliminary theorem about the second order asymptotic expansion of the log-profile penalized likelihood which prepares us for deriving the main results about the higher order structure of the penalized profile sampler in the next section.

**Theorem 1**—Let (13) and (14) be satisfied and suppose that

$$(\mathbb{P}_n - P)\ell^{(3)}(\bar{\theta}_n, \widetilde{\theta}_n, \widehat{\eta}_{\theta_n, \lambda_n}) = o_p(1), \tag{18}$$

$$\lambda_n J(\widehat{\eta}_{\theta_n, \lambda_n}) = O_p(\lambda_n + \|\widetilde{\theta}_n - \theta_0\|), \tag{19}$$

for any sequence  $\widetilde{\theta}_n$  and  $\bar{\theta}_n$  satisfying  $\widetilde{\theta}_n = \theta_0 + o_p(1)$  and  $\bar{\theta}_n = \theta_0 + o_p(1)$ . If  $\theta_0$  is an interior point in  $\Theta$  and  $\widehat{\theta}_{\lambda_n}$  is consistent, then we have

$$\sqrt{n}(\widehat{\theta}_{\lambda_n} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \widetilde{I}_0^{-1} \ell_0(X_i) + O_p(n^{1/2} \lambda_n^2), \tag{20}$$

$$\log pl_{\lambda_n}(\widetilde{\theta}_n) = \log pl_{\lambda_n}(\widehat{\theta}_{\lambda_n}) - \frac{n}{2} (\widetilde{\theta}_n - \widehat{\theta}_{\lambda_n})^T \widetilde{I}_0(\widetilde{\theta}_n - \widehat{\theta}_{\lambda_n}) + O_p(g_{\lambda_n}(\|\widetilde{\theta}_n - \widehat{\theta}_{\lambda_n}\|)), \tag{21}$$

where  $g_{\lambda_n}(w) = nw^3 + nw^2 \lambda_n + nw \lambda_n^2 + n^{1/2} \lambda_n^2$ , provided the efficient information  $\widetilde{I}_0$  is positive definite.

For the verification of (18), we need to make use of a Glivenko-Cantelli theorem for classes of functions that change with  $n$  which is a modification of theorem 2.4.3 in [22] and is explained in the appendix. Moreover, (19) implies that  $J(\widehat{\eta}_{\lambda_n}) = O_p(1)$  if the  $\widehat{\theta}_{\lambda_n}$  is asymptotically normal, which has been shown in (20).

**Remark 2**—The results in theorem 1 are useful in their own right for inference about  $\theta$ . (20) is a second higher order frequentist result in penalized semiparametric estimation regarding the asymptotic linearity of the maximum penalized likelihood estimator of  $\theta$ .

### 4 Main Results and Implications

We now state the main results on the penalized posterior profile distribution. A preliminary result, theorem 2 with corollary 1 below, shows that the penalized posterior profile distribution is asymptotically close enough to the distribution of a normal random variable with mean



$\hat{\theta}_{\lambda_n}$  and variance  $(n\tilde{I}_0)^{-1}$  with second order accuracy, which is controlled by the smoothing parameter. Similar conclusions also hold for the penalized posterior moments. Another main result, theorem 3, shows that the penalized posterior profile log-likelihood can be used to achieve second order accurate frequentist inference for  $\theta$ .

Let  $\tilde{P}_{\theta|\tilde{X}}^{\lambda_n}$  be the penalized posterior profile distribution of  $\theta$  with respect to the prior  $\rho(\theta)$ . Define

$$\Delta_{\lambda_n}(\theta) = n^{-1} \{ \log pl_{\lambda_n}(\theta) - \log pl_{\lambda_n}(\hat{\theta}_{\lambda_n}) \}.$$

**Theorem 2**

Let (20) and (21) be satisfied and suppose that

$$\Delta_{\lambda_n}(\tilde{\theta}_n) = o_p(1) \text{ implies } \tilde{\theta}_n = \theta_0 + o_p(1), \tag{22}$$

for every random  $\{\tilde{\theta}_n\} \in \Theta$ . If proper prior  $\rho(\theta_0) > 0$  and  $\rho(\cdot)$  has continuous and finite first order derivative in some neighborhood of  $\theta_0$ , then we have,

$$\sup_{\xi \in \mathbb{R}^d} \left| \tilde{P}_{\theta|\tilde{X}}^{\lambda_n} (\sqrt{n} \tilde{I}_0^{-1/2} (\theta - \hat{\theta}_{\lambda_n}) \leq \xi) - \Phi_d(\xi) \right| = O_p(n^{1/2} \lambda_n^2), \tag{23}$$

where  $\Phi_d(\cdot)$  is the distribution of the d-dimensional standard normal random variable.

**Corollary 1**

Under the assumptions of theorem 2, we have that if  $\theta$  has finite second absolute moment, then

$$\hat{\theta}_{\lambda_n} = E_{\theta|\tilde{X}}^{\lambda_n}(\theta) + O_p(\lambda_n^2), \tag{24}$$

$$\tilde{I}_0 = n^{-1} (Var_{\theta|\tilde{X}}^{\lambda_n}(\theta))^{-1} + O_p(n^{1/2} \lambda_n^2), \tag{25}$$

where  $E_{\theta|\tilde{X}}^{\lambda_n}(\theta)$  and  $Var_{\theta|\tilde{X}}^{\lambda_n}(\theta)$  are the penalized posterior profile mean and penalized posterior profile covariance matrix, respectively.

We now present another second order asymptotic frequentist property of the penalized profile sampler in terms of quantiles. The  $\alpha$ -th quantile of the penalized posterior profile distribution,

$\tau_{n\alpha}$ , is defined as  $\tau_{n\alpha} = \inf \{ \xi : \tilde{P}_{\theta|\tilde{X}}^{\lambda_n}(\theta \leq \xi) \geq \alpha \}$ , where the inf is taken componentwise. Without loss of generality, we can assume  $\tilde{P}_{\theta|\tilde{X}}^{\lambda_n}(\theta \leq \tau_{n\alpha}) = \alpha$  because of the assumed smoothness of both the prior and the likelihood in our setting. We can also define  $\kappa_{n\alpha} \equiv \sqrt{n}(\tau_{n\alpha} - \hat{\theta}_{\lambda_n})$ , i.e.,

$P_{\tilde{\theta}_X}^{\sim \lambda_n}(\sqrt{n}(\theta - \widehat{\theta}_{\lambda_n}) \leq \kappa_{n\alpha}) = \alpha$ . Note that neither  $\tau_{n\alpha}$  nor  $\kappa_{n\alpha}$  are unique if the dimension of  $\theta$  is larger than one.

**Theorem 3**

Under the assumptions of theorem 2 and assuming that  $\tilde{\ell}_0(X)$  has finite third moment with a nondegenerate distribution, then there exists a  $\widehat{\kappa}_{n\alpha}$  based on the data such that

$$P(\sqrt{n}(\widehat{\theta}_{\lambda_n} - \theta_0) \leq \widehat{\kappa}_{n\alpha}) = \alpha \text{ and } \widehat{\kappa}_{n\alpha} - \kappa_{n\alpha} = O_p(n^{1/2} \lambda_n^2) \text{ for each choice of } \kappa_{n\alpha}.$$

**Remark 3**

Theorem 3 ensures that there exists a unique  $\alpha$ -th quantile for  $\theta$  up to  $O_p(\lambda_n^2)$  in the frequentist set-up for each fixed  $\tau_{n\alpha}$ . Note that  $\tau_{n\alpha}$  is not unique if the dimension of  $\theta$  is larger than one.

**Remark 4**

Theorem 2, corollary 1 and theorem 3 above show that the penalized profile sampler generates second order asymptotic frequentist valid results in terms of distributions, moments and quantiles. Moreover, the second order accuracy of this procedure is controlled by the smoothing parameter.

**Remark 5**

Another interpretation for the role of  $\lambda_n$  in the penalized profile sampler is that we can view  $\lambda_n$  as the prior on  $J(\eta)$ , or on  $\eta$  to some extent. To see this, we can write  $lik_{\lambda_n}(\theta, \eta)$  in the following form:

$$lik_{\lambda_n}(\theta, \eta) = lik_n(\theta, \eta) \times \exp \left[ -\frac{J^2(\eta)}{2\left(\frac{1}{2n\lambda_n^2}\right)} \right]$$

This idea can be traced back to [23]. In other words, the prior on  $J(\eta)$  is a normal distribution with mean zero and variance  $(2\lambda_n^2 n)^{-1}$ . Hence it is natural to expect  $\lambda_n$  to have some effect on the convergence rate of  $\eta$ . Other possible priors on the functional parameter include Dirichlet and Gaussian processes which are more commonly used in nonparametric Bayesian methodology.

**5 Examples (Continued)**

We now illustrate verification of the assumptions in section 3.2 with the two examples that were introduced in section 2. Thus this section is a continuation of the earlier examples.

**5.1 Partly Linear Normal Model with Current Status Data**

In this section we verify the regularity conditions for the partly linear model with current status data as well as present a small simulation study to gain insight into the moderate sample size agreement with the asymptotic theory.

**5.1.1 Verification of conditions**—We will concentrate on the estimation of the regression coefficient  $\theta$ , considering the infinite dimensional parameter  $f \in \mathcal{H}_k^M$  as a nuisance parameter. The strengthened condition on  $\eta$ , together with the requirement that the density for the joint distribution  $(U, V, C)$  is strictly positive and finite, is necessary to verify the rate assumptions (27) and (28) in the below lemma 1. The score function of  $\theta, \ell_{\theta,f}$ , is given as follows:

$$\dot{\ell}_{\theta,f}(x) = uQ(x; \theta, f),$$

where

$$Q(X; \theta, f) = (1 - \Delta) \frac{\varphi(q_{\theta,f}(X))}{1 - \Phi(q_{\theta,f}(X))} - \Delta \frac{\varphi(q_{\theta,f}(X))}{\Phi(q_{\theta,f}(X))},$$

$q_{\theta,f}(x) = c - \theta u - f(v)$ , and  $\varphi$  is the density of a standard normal random variable. The least favorable direction at the true parameter value is:

$$h_0(v) = \frac{E_0(UQ^2(X; \theta, f) | V=v)}{E_0(Q^2(X; \theta, f) | V=v)},$$

where  $E_0$  is the expectation relative to the true parameters. The derivation of  $\ell_{\theta,f}$  and  $h_0(\cdot)$  is given in [3]. Thus, the least favorable submodel can be constructed as follows:

$$\ell(t, \theta, f) = \text{loglik}(t, f_t(\theta, f)), \tag{26}$$

where  $f_t(\theta, f) = f + (\theta - t)h_0$ . The concrete forms of  $\ell(t, \theta, \eta)$  and the related derivatives are given in [3] which considers a more rigid model with a known upper bound on the  $L_2$  norm of the  $k$ th derivative. The remaining assumptions are verified in the following three lemmas:

**Lemma 1:** Under the above set-up for the partly linear normal model with current status data, we then have for  $\lambda_n$  satisfying (3) and  $\tilde{\theta}_n \xrightarrow{P} \theta_0$ ,

$$\left\| \widehat{f}_{\tilde{\theta}_n, \lambda_n} - f_0 \right\|_2 = O_p(\lambda_n + \left| \tilde{\theta}_n - \theta_0 \right|), \tag{27}$$

$$\lambda_n J(\widehat{f}_{\tilde{\theta}_n, \lambda_n}) = O_p(\lambda_n + \left| \tilde{\theta}_n - \theta_0 \right|), \tag{28}$$

where  $\|\cdot\|_2$  represents the regular  $L_2$  norm. Moreover, if we also assume that  $f \in \{g: \|g\|_\infty + J(g) \leq M\}$  for some known  $M$ , then

$$\left\| \widehat{f}_{\tilde{\theta}_n} - f_0 \right\|_2 = O_p(n^{-k/(2k+1)} + \left| \tilde{\theta}_n - \theta_0 \right|), \tag{29}$$

provided condition (3) holds.

**Remark 6:** Lemma 1 implies that the convergence rate of the estimated nuisance parameter is slower than that of the regular nuisance parameter by comparing (27) and (29). This result is not surprising since the slower rate is the trade-off for the smoother nuisance parameter estimator. However, the advantage of the penalized profile sampler is that we can control the

convergence rate by assigning the smoothing parameter with different rates. To obtain the convergence rate of the non-penalized estimated nuisance parameter, we would need to assume that the Sobolev norm of the nuisance parameter has some known upper bound. Thus we can argue that the penalized method enables a relaxation of the assumptions needed for the nuisance parameter. Lemma 1 also indicates that  $\|\hat{f}_{\lambda_n} - f_0\|_2 = O_P(\lambda_n)$  and  $\|\hat{f}_n - f_0\|_2 = O_P(n^{-k/(k+2)})$ . Note that the convergence rate of the maximum penalized likelihood estimator,  $O_P(\lambda_n)$ , is deemed as the optimal rate in [23]. Similar remarks also hold for lemma 4 in semiparametric logistic regression model example below.

Lemma 1 and 4 imply that  $J(\hat{\eta}_{\lambda_n}) = O_P(1)$  and  $J(\hat{f}_{\lambda_n}) = O_P(1)$ , respectively. Thus the maximum likelihood estimators of the nuisance parameters in the two examples of this paper are consistent in the uniform norm, i.e.  $\|\hat{\eta}_{\lambda_n} - \eta_0\|_\infty = o_P(1)$  and  $\|\hat{f}_{\lambda_n} - f_0\|_\infty = o_P(1)$ , since the sequences  $\hat{\eta}_{\lambda_n}$  and  $\hat{f}_{\lambda_n}$  consist of smooth functions defined on a compact set with asymptotically bounded first-order derivatives.

**Lemma 2:** Under the above set-up for the partly linear normal model with current status data, assumptions (13), (14) and (18) are satisfied.

**Lemma 3:** Under the above set-up for the partly linear normal model with current status data, condition (22) is satisfied.

**5.1.2 Simulation study—**In this subsection, we conducted simulations for the partly linear model with two different sizes of smoothing parameter, i.e.  $\lambda_n = n^{-1/3}$  and  $\lambda_n = n^{-2/5}$ . Since we assume that  $f \in \mathcal{H}_2^M$  in the model, the above smoothing parameters satisfy (3). Our experience indicates that, in applications involving moderate sample sizes, specification of  $M$  is not needed and  $\lambda_n = n^{-1/3}$  ( $n^{-2/5}$ ) appears to work most of the time. Perhaps using cross validation to choose  $\lambda_n$  may improve the performance of the estimator in some settings, but evaluating this issue requires further study and is beyond the scope of the current paper. The contrast of the above two simulations agrees with our theoretical results that we can control the accuracy of inferences based on the penalized profile sampler by adjusting the related smoothing parameter.

We next discuss the computation of  $\hat{f}_{\theta, \lambda_n}$  in the simulations. For the special case of  $k = 2$ , we can use a cubic spline for estimating  $f$  given a fixed  $\theta$  and  $\lambda_n$ . In practice, we take a computational sieve approach suggested by Xiang and Wahba [24], which states that an estimate with the number of basis functions growing at least at the rate  $O(n^{1/5})$  can achieve the same asymptotic precision as the full space, see section 8.2 in [10] for details.

In the following, the simulations are run for various sample sizes under a Lebesgue prior. For each sample size, 200 datasets were analyzed. The regression coefficient is  $\theta = 1$  and  $f(v) = \sin(\pi v)$ . We generate  $U \sim Unif[0, 1]$ ,  $V \sim Unif[-1, 1]$  and  $C \sim Unif[0, 2]$ . For each dataset, Markov chains of length 20, 000 with a burn-in period of 5, 000 were generated using the Metropolis algorithm. The jumping density for the coefficient was normal with current iteration and variance tuned to yield an acceptance rate of 20% – 40%. The approximate variance of the estimator of  $\theta$  was computed by numerical differentiation with step size proportional to  $n^{-1/3}$  ( $n^{-2/5}$ ) for the model with smoothing parameter  $\lambda_n = n^{-1/3}$  ( $n^{-2/5}$ ) according to (21), see remark 1 in [3] for details.

Table 1 (2) in the below summarizes the simulation results for  $\theta$  with smoothing parameter  $\lambda_n = n^{-1/3}$  ( $n^{-2/5}$ ) giving the average across 200 samples of the penalized maximum likelihood estimate (PMLE), mean of the penalized profile sampler (CM), estimated standard errors based on MCMC ( $SE_M$ ), estimated standard errors based on numerical derivatives ( $SE_N$ ), boundaries for the two-sided 95% confidence interval for  $\theta$  generated by numerical differentiation and MCMC.  $L_M$  ( $L_N$ ) and  $U_M$  ( $U_N$ ) denote the lower and upper bound of the confidence interval

from the MCMC chain (numerical derivative). According to the above theoretical results, the terms  $n^{2/3}|PMLE - CM|$ ,  $n^{1/6}|SE_M - SE_N|$ ,  $n^{3/10}|SE_M - SE_N|$ ,  $n^{2/3}|L_M - L_N|$ ,  $n^{4/5}|L_M - L_N|$  and  $n^{2/3}|U_M - U_N|$ ,  $n^{4/5}|U_M - U_N|$  in Table 1 (2) are bounded in probability. And the realizations of these terms summarized in Table 1 and 2 clearly illustrate their boundedness. Furthermore, we can conclude that the penalized profile sampler with respect to different sizes of smoothing parameter can yield statistical inference with different degree of accuracy.

### 5.2 Semiparametric Logistic Regression

In the semiparametric logistic regression model, we can obtain the score function for  $\theta$  and  $\eta$  by similar analysis performed in the first example, i.e.  $\ell_{\theta,\eta}(x) = (y - F(\theta w + \eta(z)))w$  and  $A_{\theta,\eta}h_{\theta,\eta}(x) = (y - F(\theta w + \eta(z)))h_{\theta,\eta}(z)$  for  $J(h) < \infty$ , where  $A_{\eta,\eta}$  and  $h_{\theta,\eta}$  are the score operator for  $\eta$  and least favorable direction at  $(\theta, \eta)$ , respectively. And the least favorable direction at the true parameter is given in [15]:

$$h_0(z) = \frac{P_0[W \dot{F}(\theta_0 W + \eta_0(Z)) | Z=z]}{P_0[\dot{F}(\theta_0 W + \eta_0(Z)) | Z=z]},$$

where  $\dot{F}(u) = F(u)(1 - F(u))$ . The above assumptions plus the requirement that  $J(h_0) < \infty$  ensures the identifiability of the parameters. Thus the least favorable submodel can be written as:

$$\ell(t, \theta, \eta) = \text{loglik}(t, \eta_t(\theta, \eta)),$$

where  $\eta_t(\theta, \eta) = \eta + (\theta - t)h_0$ . By differentiating  $\ell(t, \theta, \eta)$  with respect to  $t$  or  $\theta$ , we obtain,

$$\begin{aligned} \dot{\ell}(t, \theta, \eta) &= (y - F(tw + \eta(z) + (\theta - t)h_0(z)))(w - h_0(z)), \\ \ddot{\ell}(t, \theta, \eta) &= -\dot{F}(tw + \eta(z) + (\theta - t)h_0(z))(w - h_0(z))^2, \\ \ell_{t,\theta}(t, \theta, \eta) &= -\dot{F}(tw + \eta(z) + (\theta - t)h_0(z))(w - h_0(z))h_0(z), \\ \ell^{(3)}(t, \theta, \eta) &= -\ddot{F}(tw + \eta(z) + (\theta - t)h_0(z))(w - h_0(z))^3, \\ \ell_{t,t,\theta}(t, \theta, \eta) &= -\ddot{F}(tw + \eta(z) + (\theta - t)h_0(z))(w - h_0(z))^2 h_0(z), \\ \ell_{t,\theta,\theta}(t, \theta, \eta) &= -\ddot{F}(tw + \eta(z) + (\theta - t)h_0(z))(w - h_0(z))h_0^2(z), \end{aligned}$$

where  $\ddot{F}(\cdot)$  is the second derivative of the function  $F(\cdot)$ . The rate assumptions will be shown in lemma 4. The remaining assumptions are verified in the last two lemmas:

**Lemma 4**—Under the above set-up for the semiparametric logistic regression model, we have for  $\lambda_n$  satisfying condition (3) and any  $\tilde{\theta}_n \xrightarrow{P} \theta_0$  that

$$\left\| \tilde{\eta}_{\theta_n, \lambda_n} - \eta_0 \right\|_2 = O_p(\lambda_n + \left\| \tilde{\theta}_n - \theta_0 \right\|), \tag{30}$$

$$\lambda_n J(\tilde{\eta}_{\theta_n, \lambda_n}) = O_p(\lambda_n + \left\| \tilde{\theta}_n - \theta_0 \right\|). \tag{31}$$

If we also assume that  $\eta \in \{g: \|g\|_\infty + J(g) \leq \tilde{M}\}$  for some known  $\tilde{M}$ , then

$$\left\| \widehat{\eta}_{\theta_n} - \eta_0 \right\|_2 = O_p(n^{-k/(2k+1)} + \left\| \widetilde{\theta}_n - \theta_0 \right\|), \quad (32)$$

provided condition (3) holds.

**Lemma 5**—Under the above set-up for the semiparametric logistic regression model, assumptions (13), (14) and (18) are satisfied.

**Lemma 6**—Under the above set-up for the semiparametric logistic regression model, condition (22) is satisfied.

## 6 Future Work

Our paper evaluates the penalized profile sampler method from the frequentist view and discusses the effect of the smoothing parameter on estimation accuracy. One potential problem of interest is to sharpen the upper bound for the convergence rate of the approximation error in this paper, like the typical second-order asymptotic results in Edgeworth expansions, see, for example [1]. A formal study about the higher order comparisons between the profile sampler procedure and fully Bayesian procedure [17], which assigns priors to both the finite dimensional parameter and the infinite dimensional nuisance parameter, is also interesting. We expect that the involvement of a suitable prior on the infinite dimensional parameter would at least not decrease the estimation accuracy of the parameter of interest.

Another worthwhile avenue of research is to develop analogs of the profile sampler and penalized profile sampler to likelihood estimation under model misspecification and to general M-estimation. Some first order results for this setting in the case where the nuisance parameter may not be root- $n$  consistent have been developed for a weighted bootstrap procedure in [11]. The studies about second order asymptotics under mild model misspecifications can provide theoretical insights into semiparametric model selection problems.

## Acknowledgments

The authors thank Dr. Joseph Kadane for several insightful discussions.

## References

1. Bentkus V, Gotze F, van Zwer WR. An Edgeworth Expansion for Symmetric Statistics. *Annals of Statistics* 1997;25:851–896.
2. Cheng G, Kosorok MR. Higher order semiparametric frequentist inference with the profile sampler. *Annals of Statistics*. 2007 In Press.
3. Cheng G, Kosorok MR. General Frequentist Properties of the Posterior Profile Distribution. *Annals of Statistics*. 2007 In Press.
4. Dalalyan A, Golubev G, Tsybakov A. A Penalized Maximum Likelihood and Semiparametric Second-Order Efficiency. *Annals of Statistics* 2006;34:169–201.
5. Good IJ, Gaskins RA. Non-parametric roughness penalties for probability densities. *Biometrika* 1971;58:255–277.
6. Huang J. Efficient estimation of the partly linear Cox model. *Annals of Statistics* 1999;27:1536–1563.
7. Kosorok, MR. *Introduction to Empirical Processes and Semiparametric Inference*. Springer; New York: 2008.
8. Kuo, HH. *Lecture Notes in Mathematics*. Berlin: Springer; 1975. *Gaussian Measure on Banach Spaces*.
9. Lee BL, Kosorok MR, Fine JP. The profile sampler. *Journal of the American Statistical Association* 2005;100:960–969.

10. Ma S, Kosorok MR. Penalized Log-likelihood Estimation for Partly Linear Transformation Models with Current Status Data. *Annals of Statistics* 2005;33:2256–2290.
11. Ma S, Kosorok MR. Robust semiparametric M-estimation and the weighted bootstrap. *Journal of Multivariate Analysis* 2005;96:190–217.
12. Ma S, Kosorok MR. Adaptive penalized M-estimation with current status data. *Annals of the Institute of Statistical Mathematics* 2006;58:511–526.
13. Mammen E, van de Geer S. Penalized quasi-likelihood estimation in partial linear models. *Annals of Statistics* 1997;25:1014–1035.
14. Murphy SA. Asymptotic Theory for the Frailty Model. *Annals of Statistics* 1995;23:182–198.
15. Murphy SA, Van der Vaart AW. Observed information in semiparametric models. *Bernoulli* 1999;5:381–412.
16. Murphy SA, Van der Vaart AW. Semiparametric mixtures in case-control studies. *Journal of Multivariate Analysis* 2001;79:1–32.
17. Shen X. Asymptotic normality in semiparametric and nonparametric posterior distributions. *Journal of the American Statistical Association* 2002;97:222–235.
18. Silverman BW. On the estimation of a probability density function by the maximum penalized likelihood method. *Annals of Statistics* 1982;10:795–810.
19. Silverman BW. Some aspects of the spline smoothing approach to nonparametric regression curve fitting (with discussion). *Journal of the Royal Statistical Society Series B* 1985;47:1–52.
20. van de Geer, S. *Empirical Processes in M-estimation*. Cambridge University Press; Cambridge: 2000.
21. van der Vaart AW. Maximum Likelihood Estimation with Partially Censored Observations. *Annals of Statistics* 1994;22:1896–1916.
22. van der Vaart, AW.; Wellner, JA. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer; New York: 1996.
23. Wahba, G. *Spline Models for Observational Data*. SIAM; Philadelphia: 1998.
24. Xiang D, Wahba G. Approximate smoothin spline methods for large data sets in the binary case. *ASA Proc of the Biometrics Section* :94–99.

## 7 Appendix

We first state classical definitions for the covering number (entropy number) and bracketing number (bracketing entropy number) for a class of functions, and then present some technical tools about the entropy calculations and increments of empirical processes which will be employed in the proofs that follow. The notations  $\gtrsim$  and  $\lesssim$  mean greater than, or smaller than, up to a universal constant.

### Definition

Let  $\mathcal{A}$  be a subset of a (pseudo-) metric space  $(\mathcal{L}, d)$  of real-valued functions. The  $\delta$ -covering number  $N(\delta, \mathcal{A}, d)$  of  $\mathcal{A}$  is the smallest  $N$  for which there exist functions  $a_1, \dots, a_N$  in  $\mathcal{L}$ , such that for each  $a \in \mathcal{A}$ ,  $d(a, a_j) \leq \delta$  for some  $j \in \{1, \dots, N\}$ . The  $\delta$ -bracketing number  $N_B(\delta, \mathcal{A}, d)$  is the smallest  $N$  for which there exist pairs of functions  $\{[a_j^L, a_j^U]\}_{j=1}^N \subset \mathcal{L}$ , with  $d(a_j^L, a_j^U) \leq \delta, j = 1, \dots, N$ , such that for each  $a \in \mathcal{A}$  there is a  $j \in \{1, \dots, N\}$  such that  $a_j^L \leq a \leq a_j^U$ . The  $\delta$ -entropy number ( $\delta$ -bracketing entropy number) is defined as  $H(\delta, \mathcal{A}, d) = \log N(\delta, \mathcal{A}, d)$  ( $H_B(\delta, \mathcal{A}, d) = \log N_B(\delta, \mathcal{A}, d)$ ).

T1. For each  $0 < C < \infty$  and  $\delta > 0$  we have

$$H_B(\delta, \{\eta: \|\eta\|_\infty \leq C, J(\eta) \leq C\}, \|\cdot\|_\infty) \lesssim \left(\frac{C}{\delta}\right)^{1/k}, \quad (33)$$



$$H(\delta, \{\eta: \|\eta\|_\infty \leq C, J(\eta) \leq C\}, \|\cdot\|_\infty) \lesssim \left(\frac{C}{\delta}\right)^{1/k}. \tag{34}$$

T2. Let  $\mathcal{F}$  be a class of measurable functions such that  $P f^2 < \delta^2$  and  $\|f\|_\infty \leq M$  for every  $f$  in  $\mathcal{F}$ . Then

$$E_p^* \|\mathbb{G}_n\|_{\mathcal{F}} \lesssim K(\delta, \mathcal{F}, L_2(P)) \left(1 + \frac{K(\delta, \mathcal{F}, L_2(P))}{\delta^2 \sqrt{n}} M\right),$$

where  $\|\mathbb{G}_n\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |G_n f|$  and  $K(\delta, \mathcal{F}, \|\cdot\|) = \int_0^\delta \sqrt{1 + H_B(\varepsilon, \mathcal{F}, \|\cdot\|)} d\varepsilon$ .

T3. Let  $\mathcal{F} = \{f_t: t \in T\}$  be a class of functions satisfying  $|f_s(x) - f_t(x)| \leq d(s, t)F(x)$  for every  $s$  and  $t$  and some fixed function  $F$ . Then, for any norm  $\|\cdot\|$ ,

$$N_B(2\varepsilon \|F\|, \mathcal{F}, \|\cdot\|) \leq N(\varepsilon, T, d).$$

T4. Let  $\mathcal{F}$  be a class of measurable functions  $f: \mathbf{D} \times \mathbf{W} \mapsto \mathbb{R}$  on a product of a finite set and an arbitrary measurable space  $(\mathbf{W}, \mathcal{W})$ . Let  $P$  be a probability measure on  $\mathbf{D} \times \mathbf{W}$  and let  $P_W$  be its marginal on  $\mathbf{W}$ . For every  $d \in \mathbf{D}$ , let  $\mathcal{F}_d$  be the set of functions  $w \mapsto f(d, w)$  as  $f$  ranges over  $\mathcal{F}$ . If every class  $\mathcal{F}_d$  is  $P_W$ -Donsker with  $\sup_{f \in \mathcal{F}_d} |P_W f(d, W)| < \infty$  for every  $d$ , then  $\mathcal{F}$  is  $P$ -Donsker.

T5. Let  $\mathcal{F}$  be a uniformly bounded class of measurable functions such that for some measurable  $f_0$ ,  $\sup_{f \in \mathcal{F}} \|f - f_0\|_\infty < \infty$ . Moreover, assume that  $H_B(\varepsilon, \mathcal{F}; L_2(P)) \leq K \varepsilon^{-\alpha}$  for some  $K < \infty$  and  $\alpha \in (0, 2)$  and for all  $\varepsilon > 0$ . Then

$$\sup_{f \in \mathcal{F}} \left[ \frac{|(\mathbb{P}_n - P)(f - f_0)|}{\|f - f_0\|_2^{1-\alpha/2} \vee n^{(\alpha-2)/12(2+\alpha)}} \right] = O_p(n^{-1/2}).$$

T6. For a probability measure  $P$ , let  $\mathcal{F}_1$  be a class of measurable functions  $f_1: \mathcal{X} \mapsto \mathbb{R}$ , and let  $\mathcal{F}_2$  denote a class of continuous nondecreasing functions  $f_2: \mathbb{R} \mapsto [0, 1]$ . Then,

$$H_B(\varepsilon, \mathcal{F}_2(\mathcal{F}_1), L_2(P)) \leq 2H_B(\varepsilon/3, \mathcal{F}_1, L_2(P)) + \sup_Q H_B(\varepsilon/3, \mathcal{F}_2, L_2(Q)).$$

T7. Let  $\mathcal{F}$  and  $\mathcal{G}$  be classes of measurable functions. Then for any probability measure  $Q$  and any  $1 \leq r \leq \infty$ ,

$$H_B(2\varepsilon, \mathcal{F} + \mathcal{G}, L_r(Q)) \leq H_B(\varepsilon, \mathcal{F}, L_r(Q)) + H_B(\varepsilon, \mathcal{G}, L_r(Q)), \tag{35}$$

and, provided  $\mathcal{F}$  and  $\mathcal{G}$  are bounded by 1 in terms of  $\|\cdot\|_\infty$ ,

$$H_B(2\varepsilon, \mathcal{F} \cdot \mathcal{G}, L_r(Q)) \leq H_B(\varepsilon, \mathcal{F}, L_r(Q)) + H_B(\varepsilon, \mathcal{G}, L_r(Q)), \tag{36}$$

where  $\mathcal{F} \cdot \mathcal{g} \equiv \{f \times g : f \in \mathcal{F} \text{ and } g \in \mathcal{g}\}$ .

**Remark 7**

The proof of T1 is found in [22]. T1 implies that the Sobolev class of functions with known bounded Sobolev norm is P-Donsker. T2 and T3 are separately lemma 3.4.2 and theorem 2.7.11 in [22]. T4 is lemma 9.2 in [16]. T5 is a result presented on page 79 of [20] and is a special case of lemma 5.13 on the same page, the proof of which can be found in pages 79–80. T6 and T7 are separately lemma 15.2 and 9.24 in [7].

**Proof of theorem 1**—We first show (20), and then we need to state one lemma before proceeding to the proof of (21). For the proof of (20), note that

$$0 = \mathbb{P}_n \ell(\widehat{\theta}_{\lambda_n}, \widehat{\theta}_{\lambda_n}, \widehat{\eta}_{\lambda_n}) + 2\lambda_n^2 \int_z \widehat{\eta}_{\lambda_n}^{(k)}(z) h_0^{(k)}(z) dz.$$

Combining the third order Taylor expansion of  $\mathbb{P}_n \ell(\widehat{\theta}_{\lambda_n}, \theta, \eta)$  around  $\theta_0$ , where  $\theta = \widehat{\theta}_{\lambda_n}$ , and  $\eta = \widehat{\eta}_{\lambda_n}$ , with conditions (13), (14) and (18), the first term in the right-hand-side of the above displayed equality equals  $\mathbb{P}_n \ell_0 - \tilde{I}_0(\widehat{\theta}_{\lambda_n} - \theta_0) + O_P(\lambda_n + \|\widehat{\theta}_{\lambda_n} - \theta_0\|)^2$ . By the inequality  $2\lambda_n^2 \int_z \widehat{\eta}_{\lambda_n}^{(k)}(z) h_0^{(k)}(z) dz \leq \lambda_n^2 (J^2(\widehat{\eta}_{\lambda_n}) + J^2(h_0))$  and assumption (19), the second term in the right-hand-side of the above equality is equal to  $O_P(\lambda_n + \|\widehat{\theta}_{\lambda_n} - \theta_0\|)^2$ . Combining everything, we obtain the following:

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{\ell}_0(X_i) = \sqrt{n}(\widehat{\theta}_{\lambda_n} - \theta_0) + O_P(n^{1/2}(\lambda_n + \|\widehat{\theta}_{\lambda_n} - \theta_0\|)^2). \tag{37}$$

The right-hand-side of (37) is of the order  $O_P(\sqrt{n}\lambda_n^2 + \sqrt{n}w_n(1+w_n+\lambda_n))$ , where  $w_n$  represents  $\|\widehat{\theta}_{\lambda_n} - \theta_0\|$ . However, its left-hand-side is trivially  $O_P(1)$ . Considering the fact that  $\sqrt{n}\lambda_n^2 = o_P(1)$ , we can deduce that  $\widehat{\theta}_{\lambda_n} - \theta_0 = O_P(n^{-1/2})$ . Inserting this into the previous display completes the proof of (20).

We next prove (21). Note that  $\widehat{\theta}_{\lambda_n} - \theta_0 = O_P(n^{-1/2})$ . Hence the order of the remainder terms in (13) and (14) become  $O_P(\lambda_n + \|\widehat{\theta}_{\lambda_n} - \theta_0\|)^2$  and  $O_P(\lambda_n + \|\widehat{\theta}_{\lambda_n} - \theta_0\|)$ , respectively. Expression (56) in lemma 7 below implies that

$$\log pl_{\lambda_n}(\widehat{\theta}_{\lambda_n}) = \log pl_{\lambda_n}(\theta_0) + n(\widehat{\theta}_{\lambda_n} - \theta_0)^T \mathbb{P}_n \tilde{\ell}_0 - \frac{n}{2}(\widehat{\theta}_{\lambda_n} - \theta_0)^T \tilde{I}_0(\widehat{\theta}_{\lambda_n} - \theta_0) + O_P(n^{1/2}\lambda_n^2). \tag{38}$$

The difference between (38) and (56) generates

$$\log pl_{\lambda_n}(\tilde{\theta}_n) = \log pl_{\lambda_n}(\widehat{\theta}_{\lambda_n}) + n(\tilde{\theta}_n - \widehat{\theta}_{\lambda_n})^T \left( \mathbb{P}_n \tilde{\ell}_0 - \tilde{I}_0(\widehat{\theta}_{\lambda_n} - \theta_0) \right) - \frac{n}{2}(\tilde{\theta}_n - \widehat{\theta}_{\lambda_n})^T \tilde{I}_0(\tilde{\theta}_n - \widehat{\theta}_{\lambda_n}) + O_P(g_{\lambda_n}(\|\tilde{\theta}_n - \widehat{\theta}_{\lambda_n}\|)).$$

(21) is now immediately obtained after considering (20).

**Proof of theorem 2**—Suppose that  $F_{\lambda_n}(\cdot)$  is the penalized posterior profile distribution of  $\sqrt{n}\varrho_n$  with respect to the prior  $\rho(\theta)$ , where the vector  $\varrho_n$  is defined as  $I_0^{-1/2}(\theta - \widehat{\theta}_{\lambda_n})$ . The parameter set for  $\varrho_n$  is  $\Xi_n$ .  $F_{\lambda_n}(\cdot)$  can be expressed as:

$$F_{\lambda_n}(\xi) = \frac{\int_{\varrho_n \in (-\infty, n^{-1/2}\xi] \cap \Xi_n} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n) \frac{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n)}{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n})} d\varrho_n}{\int_{\varrho_n \in \Xi_n} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n) \frac{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n)}{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n})} d\varrho_n}. \tag{39}$$

Note that  $d\varrho_n$  in the above is the short notation for  $d\varrho_{n1} \times \dots \times d\varrho_{nd}$ . To prove theorem 2, we first partition the parameter set  $\Xi_n$  as  $\{\Xi_n \cap \{\|\varrho_n\| > r_n\}\} \cup \{\Xi_n \cap \{\|\varrho_n\| \leq r_n\}\}$ . By choosing the proper order of  $r_n$ , we find the posterior mass in the first partition region is of arbitrarily small order, as verified in lemma 2.1 immediately below, and the mass inside the second partition region can be approximated by a stochastic polynomial in powers of  $n^{-1/2}$  with error of order dependent on the smoothing parameter, as verified in lemma 2.2 below. This basic technique applies to both the denominator and the numerator, yielding the quotient series, which gives the desired result.

**lemma 2.1**—Choose  $r_n = o(n^{-1/3})$  and  $\sqrt{n}r_n \rightarrow \infty$ . Under the conditions of theorem 2, we have

$$\int_{\|\varrho_n\| > r_n} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n) \frac{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n)}{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n})} d\varrho_n = O_p(n^{-M}), \tag{40}$$

for any positive number  $M$ .

**Proof**—Fix  $r > 0$ . We then have

$$\int_{\|\varrho_n\| > r} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n) \frac{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-1/2}\varrho_n)}{pl_{\lambda_n}(\widehat{\theta}_{\lambda_n})} d\varrho_n \leq I\{\Delta_{\lambda_n}^r < -n^{-1/2}\} \exp(-\sqrt{n}) \int_{\Theta} \rho(\theta) d\theta + I\{\Delta_{\lambda_n}^r \geq -n^{-1/2}\},$$

where  $\Delta_{\lambda_n}^r = \sup_{\|\varrho_n\| > r} \Delta_{\lambda_n}(\widehat{\theta}_{\lambda_n} + \varrho_n I_0^{-1/2})$ . According to lemma 3.2 in [2],

$I\{\Delta_{\lambda_n}^r \geq -n^{-1/2}\} = O_p(n^{-M})$  for any positive decreasing  $r \rightarrow 0$ . Note that the above inequality holds uniformly for any decreasing  $r_n \rightarrow 0$ . Therefore, we can choose a positive decreasing sequence  $r_n = o(n^{-1/3})$  with  $\sqrt{n}r_n \rightarrow \infty$  such that (40) holds.

**lemma 2.2**—Choose  $r_n = o(n^{-1/3})$  and  $\sqrt{n}r_n \rightarrow \infty$ . Under the conditions of theorem 2, we have

$$\int_{\|\varrho_n\| \leq r_n} \left| \frac{p_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} \varrho_n)}{p_{\lambda_n}(\widehat{\theta})} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} \varrho_n) - \exp\left(-\frac{n}{2} \varrho_n^T \varrho_n\right) \rho(\widehat{\theta}_{\lambda_n}) \right| \times d\varrho_n = O_p(n^{-(d-1)/2} \lambda_n^2). \tag{41}$$

**Proof**—The posterior mass over the region  $\|\varrho_n\|_2 \leq r_n$  is bounded by

$$\begin{aligned} & \int_{\|\varrho_n\|_2 \leq r_n} \left| \frac{p_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} \varrho_n)}{p_{\lambda_n}(\widehat{\theta}_{\lambda_n})} \rho(\widehat{\theta}_{\lambda_n}) - \exp\left(-\frac{n}{2} \varrho_n^T \varrho_n\right) \rho(\widehat{\theta}_{\lambda_n}) \right| d\varrho_n \tag{*} \\ & + \int_{\|\varrho_n\|_2 \leq r_n} \left| \frac{p_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} \varrho_n)}{p_{\lambda_n}(\widehat{\theta}_{\lambda_n})} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} \varrho_n) - \frac{p_{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} \varrho_n)}{p_{\lambda_n}(\widehat{\theta}_{\lambda_n})} \rho(\widehat{\theta}_{\lambda_n}) \right| d\varrho_n. \tag{**} \end{aligned}$$

By (21), we obtain

$$\tag{*} = \int_{\|\varrho_n\|_2 \leq r_n} \left[ \rho(\widehat{\theta}_{\lambda_n}) \exp\left(-\frac{n \varrho_n^T \varrho_n}{2}\right) |\exp(O_p(g_{\lambda_n}(\|\varrho_n\|))) - 1| \right] d\varrho_n.$$

Obviously the order of (\*) depends on that of  $|\exp(O_p(g_{\lambda_n}(\|\varrho_n\|))) - 1|$  for  $\lambda_n$  satisfying (3) and  $\|\varrho_n\| \leq r_n$ . In order to analyze its order, we partition the set  $\{\lambda_n = o_p(n^{-1/4})$  and  $\lambda_n^{-1} = O_p(n^{k/(2k+1)})\}$  with the set  $\{\lambda_n = O_p(n^{-1/3})\}$ , i.e.

$U_n = \{\lambda_n = o_p(n^{-1/4})$  and  $\lambda_n^{-1} = O_p(n^{k/(2k+1)})\} \cap \{\lambda_n = O_p(n^{-1/3})\}$  and

$L_n = \{\lambda_n = o_p(n^{-1/4})$  and  $\lambda_n^{-1} = O_p(n^{k/(2k+1)})\} \cap \{\lambda_n = O_p(n^{-1/3})\}^C$ . For the set  $U_n$ , we have  $|\exp(O_p(g_{\lambda_n}(\|\varrho_n\|))) - 1| = g_{\lambda_n}(\|\varrho_n\|) \times O_p(1)$ . For the set  $L_n$ , we have

$O_p(g_{\lambda_n}(\|\varrho_n\|)) = O_p(n \|\varrho_n\| \lambda_n^2 + n^{1/2} \lambda_n^2)$ . We can take  $r_n = n^{-1-\delta} \lambda_n^{-2}$  for some  $\delta > 0$  such that  $\sqrt{n} r_n \rightarrow \infty$  and  $r_n = o(n^{-1/3})$ . Then  $|\exp(O_p(g_{\lambda_n}(\|\varrho_n\|))) - 1| = (n \|\varrho_n\| \lambda_n^2 + n^{1/2} \lambda_n^2) \times O_p(1)$ .

Combining with the above, we know that  $\tag{*} = O_p(n^{-(d-1)/2} \lambda_n^2)$ . By similar analysis, we can also show that  $\tag{**}$  has the same order. This completes the proof of lemma 2.2.

We next start the formal proof of theorem 2. By considering both lemma 2.1 and lemma 2.2, we know the denominator of (39) equals

$$\int_{\{\|\varrho_n\|_2 \leq r_n\} \cap \Xi_n} \left[ \exp\left(-\frac{n}{2} \varrho_n^T \varrho_n\right) \rho(\widehat{\theta}_{\lambda_n}) \right] d\varrho_n + O_p(n^{-(d-1)/2} \lambda_n^2).$$

The first term in the above display equals

$$n^{-d/2} \rho(\widehat{\theta}_{\lambda_n}) \int_{\{\|u_n\|_2 \leq \sqrt{n} r_n\} \cap \sqrt{n} \Xi_n} e^{-u_n^T u_n / 2} du_n = n^{-d/2} \rho(\widehat{\theta}_{\lambda_n}) \int_{\mathbb{R}^d} e^{-u_n^T u_n / 2} du_n + O(n^{-(d-1)/2} \lambda_n^2),$$

where  $u_n = \sqrt{n}Q_n$ . The above equality follows from the inequality that  $\int_x^\infty e^{-y^2/2} dy \leq x^{-1} e^{-x^2/2}$  for any  $x > 0$ . Consolidating the above analyses, we deduce that the denominator of (39) equals  $n^{-\frac{d}{2}} \rho(\widehat{\theta}_{\lambda_n})(2\pi)^{d/2} + O_p(n^{-(d-1)/2} \lambda_n^2)$ . The same analysis also applies to the numerator, thus completing the whole proof.

**Proof of corollary 1**—We only show (24) in what follows. (25) can be verified similarly.

Showing (24) is equivalent to establishing  $\widetilde{E}_{\theta,x}^{\lambda_n}(Q_n) = O_p(\lambda_n^2)$ . Note that  $\widetilde{E}_{\theta,x}^{\lambda_n}(Q_n)$  can be written as:

$$\widetilde{E}_{\theta,x}^{\lambda_n}(Q_n) = \frac{\int_{Q_n \in \Xi_n} Q_n \rho(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} Q_n) \frac{\rho^{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} Q_n)}{\rho^{\lambda_n}(\widehat{\theta}_{\lambda_n})} dQ_n}{\int_{Q_n \in \Xi_n} \rho(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} Q_n) \frac{\rho^{\lambda_n}(\widehat{\theta}_{\lambda_n} + I_0^{-\frac{1}{2}} Q_n)}{\rho^{\lambda_n}(\widehat{\theta}_{\lambda_n})} dQ_n}.$$

By analysis similar to that applied in the proof of theorem 2, we know the denominator in the above display is  $n^{-d/2} (2\pi)^{d/2} \rho(\widehat{\theta}_{\lambda_n}) + O_p(n^{-(d-1)/2} \lambda_n^2)$  and the numerator is a random vector of order  $O_p(n^{-d/2} \lambda_n^2)$ . This yields the conclusion.

**Proof of theorem 3**—Note that (23) implies  $\kappa_{n\alpha} = I_0^{-1/2} z_\alpha + O_p(n^{1/2} \lambda_n^2)$ , for any  $\xi < \alpha < 1 - \xi$ , where  $\xi \in (0, \frac{1}{2})$ . Note also that the  $\alpha$ -th quantile of a  $d$  dimensional standard normal distribution,  $z_\alpha$ , is not unique if  $d > 1$ . The classical Edgeworth expansion implies that

$P(n^{-1/2} \sum_{i=1}^n I_0^{-1/2} \ell_0(X_i) \leq z_\alpha + a_n(\alpha)) = \alpha$ , where  $a_n(\alpha) = O(n^{-1/2})$ , for  $\xi < \alpha < 1 - \xi$ . Note that  $a_n(\alpha)$  is uniquely determined for each fixed  $z_\alpha$  since  $\ell_0(X_i)$  has at least one absolutely continuous component. Let  $\widehat{\kappa}_{n\alpha} = I_0^{-1/2} z_\alpha + (\sqrt{n}(\widehat{\theta}_{\lambda_n} - \theta_0) - n^{-1/2} \sum_{i=1}^n I_0^{-1/2} \ell_0(X_i)) + I_0^{-1/2} a_n(\alpha)$ . Then  $P(\sqrt{n}(\widehat{\theta}_{\lambda_n} - \theta_0) \leq \widehat{\kappa}_{n\alpha}) = \alpha$ . Combining with (20), we obtain  $\widehat{\kappa}_{n\alpha} = \kappa_{n\alpha} + O_p(n^{1/2} \lambda_n^2)$ . The uniqueness of  $\widehat{\kappa}_{n\alpha}$  up to order  $O_p(n^{1/2} \lambda_n^2)$  follows from that of  $a_n(\alpha)$  for each chosen  $z_\alpha$ .

**Proof of lemma 1**—We first present a technical lemma before the formal proof of lemma 1. In lemma 1.1 we define

$$\mathcal{K} = \left\{ \frac{\ell_{\theta,\eta}(X) - \ell_0(X)}{1+J(\eta)} : \|\theta - \theta_0\| \leq C_1, \|\eta - \eta_0\|_\infty \leq C_1, J(\eta) < \infty \right\},$$

for a known constant  $C_1 < \infty$ . Combining with T5, we use condition (42) below to control the order of the increments of the empirical processes indexed by  $\ell_{\theta,\eta}$ :

$$H_B(\varepsilon, \mathcal{K}, L_2(P)) \lesssim \varepsilon^{-1/k}. \tag{42}$$

We next assume two smoothness conditions about the criterion function  $(\theta, \eta) \mapsto P\ell_{\theta,\eta}$ , i.e.,

$$\|\ell_{\theta,\eta} - \ell_0\|_2 \lesssim \|\theta - \theta_0\| + d_\theta(\eta, \eta_0), \tag{43}$$

$$P(\ell_{\theta,\eta} - \ell_{\theta,\eta_0}) \lesssim -d_{\theta}^2(\eta, \eta_0) + \|\theta - \theta_0\|^2. \tag{44}$$

Here  $d_{\theta}^2(\eta, \eta_0)$  can be thought of as the square of a distance, but the following lemma is valid for arbitrary functions  $\eta \mapsto d_{\theta}^2(\eta, \eta_0)$ . Finally, we assume a somewhat stronger assumption on the density, i.e.,

$$p_{\theta,\eta}/p_{\theta,\eta_0} \text{ is bounded away from zero and infinity.} \tag{45}$$

But (45) is trivial to satisfy in our first model.

**Lemma 1.1**—Assume conditions (42)–(45) in the above hold for every  $\theta \in \Theta_n$  and  $\eta \in \mathcal{V}_n$ . Then we have

$$\begin{aligned} d_{\theta_n}(\tilde{\eta}_{\theta_n, \lambda_n}, \eta_0) &= O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|), \\ \lambda_n J(\tilde{\eta}_{\theta_n, \lambda_n}) &= O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|), \end{aligned}$$

for  $(\tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})$  satisfying  $P(\tilde{\theta}_n \in \Theta_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n} \in \mathcal{V}_n) \rightarrow 1$ .

**Proof of lemma 1.1**—The definition of  $\hat{\eta}_{\tilde{\theta}_n, \lambda_n}$  implies that

$$\begin{aligned} \lambda_n^2 J^2(\tilde{\eta}_{\theta_n, \lambda_n}) &\leq \lambda_n^2 J^2(\eta_0) + (\mathbb{P}_n - P) \left( \ell_{\theta_n, \tilde{\eta}_{\theta_n, \lambda_n}} - \ell_{\theta_n, \eta_0} \right) + P \left( \ell_{\theta_n, \tilde{\eta}_{\theta_n, \lambda_n}} - \ell_{\theta_n, \eta_0} \right) \\ &\leq \lambda_n^2 J^2(\eta_0) + I + II. \end{aligned}$$

Note that by T5 and assumption (42), we have

$$\begin{aligned} I &\leq (1 + J(\tilde{\eta}_{\theta_n, \lambda_n})) O_p(n^{-1/2}) \times \left\{ \left\| \frac{\ell_{\theta_n, \tilde{\eta}_{\theta_n, \lambda_n}} - \ell_{\theta_n, \eta_0}}{1 + J(\tilde{\eta}_{\theta_n, \lambda_n})} \right\|_2^{1 - \frac{1}{2k}} \vee n^{-\frac{2k-1}{2(2k+1)}} \right\} \\ &\quad + (1 + J(\eta_0)) O_p(n^{-1/2}) \times \left\{ \left\| \frac{\ell_{\theta_n, \eta_0} - \ell_{\theta_n, \eta_0}}{1 + J(\eta_0)} \right\|_2^{1 - \frac{1}{2k}} \vee n^{-\frac{2k-1}{2(2k+1)}} \right\}. \end{aligned}$$

By assumption (44), we have

$$II \lesssim -d_{\theta_n}^2(\tilde{\eta}_{\theta_n, \lambda_n}, \eta_0) + \|\tilde{\theta}_n - \theta_0\|^2.$$

Combining with the above, we can deduce that

$$\begin{aligned} \widehat{d}_n^2 + \lambda_n^2 \widehat{J}_n^2 &\lesssim (1 + \widehat{J}_n) O_p(n^{-1/2}) \times \left\{ \left( \frac{\widehat{d}_n + \|\tilde{\theta}_n - \theta_0\|}{1 + \widehat{J}_n} \right)^{1 - \frac{1}{2k}} \vee n^{-\frac{2k-1}{2(2k+1)}} \right\} \\ &+ (1 + J_0) O_p(n^{-1/2}) \times \left\{ \left( \frac{\|\tilde{\theta}_n - \theta_0\|}{1 + J_0} \right)^{1 - \frac{1}{2k}} \vee n^{-\frac{2k-1}{2(2k+1)}} \right\} \\ &+ \lambda_n^2 J_0^2 + \|\tilde{\theta}_n - \theta_0\|^2, \end{aligned} \tag{46}$$

where  $\widehat{d}_n = d_{\tilde{\theta}_n}(\widehat{\eta}_{\tilde{\theta}_n, \lambda_n}, \eta_0)$ ,  $J(\eta_0) = J_0$  and  $\widehat{J}_n = J(\widehat{\eta}_{\tilde{\theta}_n, \lambda_n})$ . The above inequality follows from assumption (43). Combining all of the above inequalities, we can deduce that

$$u_n^2 = O_p(1) + O_p(1) u_n^{1 - \frac{1}{2k}}, \tag{47}$$

$$v_n = v_n^{-1} O_p(\|\tilde{\theta}_n - \theta_0\|^2) + u_n^{1 - \frac{1}{2k}} O_p(\lambda_n) + O_p(n^{-\frac{1}{2}} \lambda_n^{-1} \|\tilde{\theta}_n - \theta_0\|^{1 - \frac{1}{2k}}), \tag{48}$$

where  $u_n = (\widehat{d}_n + \|\tilde{\theta}_n - \theta_0\|)/(\lambda_n + \lambda_n \widehat{J}_n)$  and  $v_n = \lambda_n \widehat{J}_n + \lambda_n$ . The equation (47) implies that  $u_n = O_p(1)$ . Inserting  $u_n = O_p(1)$  into (48), we can know that  $v_n = O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)$ , which implies  $u_n$  has the desired order. This completes the whole proof.

We now apply lemma 1.1 to derive the related convergence rates in the partly linear model. Conditions (43)–(45) can be verified easily in this example because  $\ell_{\theta, f}$  has finite second moment, and  $p_{\theta, f}$  is bounded away from zero and infinity uniformly for  $(\theta, f)$  ranging over the whole parameter space. Note that  $d_{\theta, f}(f_0) = \|p_{\theta, f} - p_0\|_2 \gtrsim \|q_{\theta, f} - q_{\theta_0, f_0}\|_2$  by Taylor expansion. Then by the assumption that  $PVar(U|V)$  is positive definite, we know that  $\|q_{\tilde{\theta}_n, \widehat{\eta}_{\tilde{\theta}_n, \lambda_n}} - q_{\theta_0, f_0}\|_2 = O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)$  implies  $\|\widehat{f}_{\tilde{\theta}_n, \lambda_n} - f_0\|_2 = O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)$ . Thus we only need to show that the  $\varepsilon$ -bracketing entropy number of the function class  $\mathcal{O}$  defined below is of order  $\varepsilon^{-1/k}$  to complete the proof of (27)–(28):

$$\mathcal{O} \equiv \left\{ \frac{\ell_{\theta, f}(X)}{1 + J(f)} : \|\theta - \theta_0\| \leq C_1, \|f - f_0\|_\infty \leq C_1, J(f) < \infty \right\},$$

for some constant  $C_1$ . Note that  $\ell_{\theta, f}(X)/(1 + J(f))$  can be rewritten as:

$$\Delta A^{-1} \log \Phi(\bar{q}_{\theta, f} A) + (1 - \Delta) A^{-1} \log(1 - \Phi(\bar{q}_{\theta, f} A)), \tag{49}$$

where  $A = 1 + J(f)$  and  $\bar{q}_{\theta, f} \in \mathcal{O}_1$ , where

$$\mathcal{O}_1 \equiv \left\{ \frac{q_{\theta, f}(X)}{1 + J(f)} : \|\theta - \theta_0\| \leq C_1, \|f - f_0\|_\infty \leq C_1, J(f) < \infty \right\},$$

and where we know  $H_B(\varepsilon, \mathcal{O}_1, L_2(P)) \lesssim \varepsilon^{-1/k}$  by T1.



We next calculate the  $\varepsilon$ -bracketing entropy number with  $L_2$  norm for the class of functions  $R_1 \equiv \{k_a(t): t \mapsto a^{-1} \log \Phi(at) \text{ for } a \geq 1 \text{ and } t \in \mathbb{R}\}$ . By some analysis we know that  $k_a(t)$  is strictly decreasing in  $a$  for  $t \in \mathbb{R}$ , and  $\sup_{t \in \mathbb{R}} |k_a(t) - k_b(t)| \lesssim |a - b|$  because  $|\partial/\partial a(k_a(t))|$  is bounded uniformly over  $t \in \mathbb{R}$ . In addition, we know that  $\sup_{a,b \geq \lambda_0, t \in \mathbb{R}} |k_a(t) - k_b(t)| \lesssim A_0^{-1}$  because the function  $u \mapsto u \log \Phi(u^{-1}t)$  has bounded derivative for  $0 < u \leq 1$  uniformly over  $t \in \mathbb{R}$ . The above two inequalities imply that the  $\varepsilon$ -bracketing number with uniform norm is of order  $O(\varepsilon^{-2})$  for  $a \in [1, \varepsilon^{-1}]$  and is 1 for  $a > \varepsilon^{-1}$ . Thus we know  $H_B(\varepsilon, R_1, L_2) = O(\log \varepsilon^{-2})$ . By applying a similar analysis to  $R_2 \equiv \{k_a(t): t \mapsto a^{-1} \log(1 - \Phi(at)) \text{ for } a \geq 1 \text{ and } t \in \mathbb{R}\}$ , we obtain that  $H_B(\varepsilon, R_2, L_2) = O(\log \varepsilon^{-2})$ . Combining this with T6 and T7, we deduce that  $H_B(\varepsilon, \mathcal{O}, L_2) \lesssim \varepsilon^{-1/k}$ . This completes the proof of (27)–(28).

For the proof of (29), we apply arguments similar to those used in the proof of lemma 1.1 but after setting  $\lambda_n, J_0$  and  $\tilde{J}_n$  to zero in (46). Then we obtain the following equality:

$$\widehat{d}_n^2 = O_p(n^{-2k/(2k+1)}) + \|\tilde{\theta}_n - \theta_0\|^2 + O_p(n^{-1/2}) \|\tilde{\theta}_n - \theta_0\|^{1-1/2k} + O_p(n^{-1/2}) (\|\tilde{\theta}_n - \theta_0\| + \widehat{d}_n)^{1-1/2k}. \text{ By treating } \|\tilde{\theta}_n - \theta_0\| \leq n^{-k/(2k+1)} \text{ and } \|\tilde{\theta}_n - \theta_0\| > n^{-k/(2k+1)} \text{ differently in the above equality, we obtain (29).}$$

**Proof of lemma 2**—Based on the discussions of (13) and (14), we need to verify the smoothness conditions and asymptotic equicontinuity conditions, i.e. (15)–(17), for the function  $\ell(t, \theta, \eta)$  and its related derivatives. The first set of conditions are verified in lemma 5 of [3]. For the verifications of (15)–(17), we first show condition (17). Without loss of generality, we assume that  $\lambda_n$  is bounded below by a multiple of  $n^{-k/(2k+1)}$  and bounded above by  $n^{-1/4}$  in view of (3). Thus

$$P \left( \frac{\dot{\ell}(\theta_0, \theta_0, \widehat{f}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}_0}{n^{\frac{1}{4k+2}}(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)} \right)^2 \lesssim \frac{\|\widehat{f}_{\tilde{\theta}_n, \lambda_n} - f_0\|_2^2}{n^{\frac{1}{2k+1}}(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)^2} = O_p \left( n^{-\frac{1}{2k+1}} \right),$$

where (27) implies the equality in the above expression.

By (28), we know that  $J(\widehat{f}_{\tilde{\theta}_n, \lambda_n}) = O_p(1 + \|\tilde{\theta}_n - \theta_0\|/\lambda_n)$  and  $\|\widehat{f}_{\tilde{\theta}_n, \lambda_n}\|_\infty$  is bounded by some constant, since  $f \in \mathcal{H}_k^M$ . We then define the set  $\mathcal{Q}_n$  as follows:

$$\left\{ \frac{\dot{\ell}(\theta_0, \theta_0, f) - \dot{\ell}_0}{n^{\frac{1}{4k+2}}(\lambda_n + \|\theta - \theta_0\|)}; J(f) \leq C_n(1 + \frac{\|\theta - \theta_0\|}{\lambda_n}), \|f\|_\infty \leq M, \|\theta - \theta_0\| \leq \delta \right\} \cap \left\{ g \in L_2(P): Pg^2 \leq C_n n^{-\frac{1}{2k+1}} \right\},$$

for some  $\delta > 0$ . Obviously the function  $n^{-1/(4k+2)}(\ell(\theta_0, \theta_0, \widehat{f}_{\tilde{\theta}_n, \lambda_n}) - \ell)/(\lambda_n + \|\tilde{\theta}_n - \theta_0\|) \in \mathcal{Q}_n$  on a set of probability arbitrarily close to one, as  $C_n \rightarrow \infty$ . If we can show  $\lim_{n \rightarrow \infty} E^* \|\mathbb{G}_n\|_{\mathcal{Q}_n} < \infty$  by T2, then assumption (17) is verified. Note that  $\dot{\ell}(\theta_0, \theta_0, f)$  depends on  $f$  in a Lipschitz manner. Consequently we can bound  $H_B(\varepsilon, \mathcal{Q}_n, L_2(P))$  by the product of some constant and  $H(\varepsilon, \mathcal{R}_n, L_2(P))$  in view of T3.  $\mathcal{R}_n$  is defined as

$$\{H_n(f): J(H_n(f)) \lesssim \lambda_n^{-1} n^{-1/(4k+2)}, \|H_n(f)\|_\infty \lesssim \lambda_n^{-1} n^{-1/(4k+2)}\},$$

where  $H_n(f) = f/(n^{1/(4k+2)}(\lambda_n + \|\theta - \theta_0\|))$ . By [22],

we know that

$$H(\varepsilon, \mathcal{R}_n, L_2(P)) \lesssim (\lambda_n^{-1} n^{-\frac{1}{4k+2}}) / \varepsilon)^{1/k}.$$

Note that  $\delta_n = n^{-1/(4k+2)}$  and  $M_n = n^{(2k-1)/(4k+2)}$  in T2. Thus by calculation we know that  $K(\delta_n, \mathcal{Q}_n, L_2(P)) \lesssim \lambda_n^{-1/2k} n^{-1/(4k+2)}$ . Then by T2 we can show that  $\lim_{n \rightarrow \infty} E^* \|\mathbb{G}_n\|_{\mathcal{Q}_n} < \infty$ .

For the proof of (15), we only need to show (15) holds for  $\tilde{\theta}_n = \hat{\theta}_n + o(n^{-1/3})$  based on the arguments in lemma 2.2. We then show that

$$\mathbb{G}_n(\check{\ell}(\theta_0, \tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}) - \check{\ell}_0) = o_p(1 + n^{1/3} \|\tilde{\theta}_n - \theta_0\|) = o_p(1).$$

By the rate assumptions (27), we have

$$P \left( \frac{\check{\ell}(\theta_0, \tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}) - \check{\ell}_0}{1 + n^{1/3} \|\tilde{\theta}_n - \theta_0\|} \right)^2 \leq \frac{\|\tilde{\theta}_n - \theta_0\|^2 + \|\hat{f}_{\tilde{\theta}_n, \lambda_n} - f_0\|_2^2}{(1 + n^{1/3} \|\tilde{\theta}_n - \theta_0\|)^2} = O_p(n^{-1/2}).$$

We next define  $\mathcal{Q}_n$  as follows:

$$\left\{ \frac{\check{\ell}(\theta_0, \theta, f) - \check{\ell}_0}{1 + n^{1/3} \|\theta - \theta_0\|} : J(f) \leq C_n \left( 1 + \frac{\|\theta - \theta_0\|}{\lambda_n} \right), \|f\|_\infty \leq M, \|\theta - \theta_0\| < \delta \right\} \cap \left\{ g \in L_2(P) : P g^2 \leq C_n n^{-\frac{1}{2}} \right\}.$$

Obviously the function  $(\check{\ell}(\theta_0, \tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}) - \check{\ell}_0) / (1 + n^{1/3} \|\tilde{\theta}_n - \theta_0\|) \in \mathcal{Q}_n$  on a set of probability arbitrarily close to one, as  $C_n \rightarrow \infty$ . If we can show  $\lim_{n \rightarrow \infty} E^* \|\mathbb{G}_n\|_{\mathcal{Q}_n} \rightarrow 0$  by T2, then the proof of (15) is completed. Accordingly, note that  $\check{\ell}(\theta_0, \theta, f)$  depends on  $(\theta, f)$  in a Lipschitz manner. Consequently we can bound  $H_B(\varepsilon, \mathcal{Q}_n, L_2(P))$  by the product of some constant and  $(H(\varepsilon, \overline{\mathcal{R}}_n, L_2(P)) + \log(1/\varepsilon))$  in view of T3.  $\overline{\mathcal{R}}_n$  is defined as

$$\{H_n(f) : J(H_n(f)) \lesssim 1 + (n^{1/3} \lambda_n)^{-1}, \|H_n(f)\|_\infty \lesssim 1 + (n^{1/3} \lambda_n)^{-1}\},$$

where  $H_n(f) = f / (1 + n^{1/3} \|\theta - \theta_0\|)$ . By [22], we know that

$$H(\varepsilon, \overline{\mathcal{R}}_n, L_2(P)) \lesssim ((1 + n^{-1/3} \lambda_n^{-1}) / \varepsilon)^{1/k}.$$

Then by analysis similar to that used in the proof of (17), we can show that  $\lim_{n \rightarrow \infty} E^* \|\mathbb{G}_n\|_{\mathcal{Q}_n} \rightarrow 0$  in view of T2. This completes the proof of (15).

For the proof of (16), it suffices to show that  $\mathbb{G}_n(\ell_{t,\theta}(\theta_0, \tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}) - \ell_{t,\theta}(\theta_0, \theta_0, f_0)) = o_p(1)$  for  $\tilde{\theta}_n = \hat{\theta}_n + o(n^{-1/3})$  and for  $\tilde{\theta}_n$  between  $\hat{\theta}_n$  and  $\theta_0$ , in view of lemma 2.2. Then we can show that

$G_n(\ell_{t,\theta}(\theta_0, \tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}) - \ell_{t,\theta}(\theta_0, \theta_0, f_0)) = o_P(1 + n^{1/3} \|\tilde{\theta}_n - \theta_0\|) = o_P(1)$  by similar analysis as used in the proof of (15).

In the last part, we show (18). It suffices to verify that the sequence of classes of functions  $\mathcal{V}_n$  is  $P$ -Glivenko-Cantelli, where  $\mathcal{V}_n \equiv \{\ell^{(3)}(\tilde{\theta}_n, \theta_n, \hat{f}_{\tilde{\theta}_n, \lambda_n})(x)\}$ , for every random sequence  $\tilde{\theta}_n \rightarrow \theta_0$  and  $\theta_n \rightarrow \theta_0$  in probability. A Glivenko-Cantelli theorem for classes of functions that change with  $n$  is needed. By revising theorem 2.4.3 in [22] with minor notational changes, we obtain the following suitable extension of the uniform entropy Glivenko-Cantelli theorem: Let  $\tilde{\mathcal{F}}_n$  be suitably measurable classes of functions with uniformly integrable functions and  $H(\varepsilon, \tilde{\mathcal{F}}_n, L_1(\mathbb{P}_n)) = o_p^*(n)$  for any  $\varepsilon > 0$ . Then  $\|\mathbb{P}_n - P\|_{\tilde{\mathcal{F}}_n} \rightarrow 0$  in probability for every  $\varepsilon > 0$ . We then apply this revised theorem to the set  $\tilde{\mathcal{F}}_n$  of functions  $\ell^{(3)}(t, \theta, f)$  with  $t$  and  $\theta$  ranging over a neighborhood of  $\theta_0$  and  $\lambda_n J(f)$  bounded by a constant. By the form of  $\ell^{(3)}(t, \theta, f)$ , the entropy number for  $\mathcal{V}_n$  is equal to that of

$$\tilde{\mathcal{F}}_n \equiv \{\varphi(q_{t,f_i(\theta,f)}(x))R(q_{t,f_i(\theta,f)}(x)): (t, \theta) \in V_{\theta_0}, \lambda_n J(f) \leq C, \|f\|_{\infty} \leq M\}.$$

By arguments similar to those used in lemma 7.2 of [15], we know that

$\sup_Q H(\varepsilon, \tilde{\mathcal{F}}_n, L_1(Q)) \lesssim (1 + \lambda_n^{-1} / \varepsilon)^{1/k} = o_p(n)$ . Moreover, the  $\tilde{\mathcal{F}}_n$  are uniformly bounded since  $f \in \mathcal{H}_k^M$ . Considering the fact that the probability that  $\mathcal{V}_n$  is contained in  $\tilde{\mathcal{F}}_n$  tends to 1, we have completed the proof of (18).

**Proof of lemma 3**—By the assumption that  $\Delta_{\lambda_n}(\tilde{\theta}_n) = o_P(1)$ , we have  $\Delta_{\lambda_n}(\tilde{\theta}_n) - \Delta_{\lambda_n}(\theta_0) \geq o_P(1)$ . Thus the following inequality holds:

$$n^{-1} \sum_{i=1}^n \log \left[ \frac{\text{lik}(\tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}, X_i)}{\text{lik}(\theta_0, \hat{f}_{\theta_0, \lambda_n}, X_i)} \right] - \lambda_n^2 [J^2(\hat{f}_{\tilde{\theta}_n, \lambda_n}) - J^2(\hat{f}_{\theta_0, \lambda_n})] \geq o_p(1)$$

By considering assumption (19), the above inequality simplifies to

$$n^{-1} \sum_{i=1}^n \log \left[ \frac{H(\tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}; X_i)}{H(\theta_0, \hat{f}_{\theta_0, \lambda_n}; X_i)} \right] \geq o_p(1),$$

where  $H(\theta, f; X) = \Delta\Phi(C - \theta U - f(V)) + (1 - \Delta)(1 - \Phi(C - \theta U - f(V)))$ . By arguments similar to those used in lemma 2 and by T4, we know  $H(\tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}; X_i)$  belongs to some  $P$ -Donsker class. Combining the above conclusion and the inequality  $\alpha \log x \leq \log(1 + \alpha\{x - 1\})$  for some  $\alpha \in (0, 1)$  and any  $x > 0$ , we can show that

$$P \log \left[ 1 + \alpha \left( \frac{H(\tilde{\theta}_n, \hat{f}_{\tilde{\theta}_n, \lambda_n}; X_i)}{H(\theta_0, \hat{f}_{\theta_0, \lambda_n}; X_i)} - 1 \right) \right] \geq o_p(1). \tag{50}$$

The remainder of the proof follows the proof of lemma 6 in [3].

**Proof of lemma 4**—The boundedness condition (45) in Lemma 1.1 can not be satisfied in semiparametric logistic regression model. Hence we propose lemma 4.1 below to relax this condition by choosing the criterion function  $m_{\theta,\eta} = \log[(p_{\theta,\eta} + p_{\theta,\eta_0})/2p_{\theta,\eta_0}]$ . Obviously,  $m_{\theta,\eta}$  is trivially bounded away from zero. It is also bounded above for  $(\theta, \eta)$  around their true values if  $p_{\theta,\eta_0}(x)$  is bounded away from zero uniformly in  $x$  and  $p_{\theta,\eta}$  is bounded above. The first condition is satisfied if the map  $\theta \mapsto p_{\theta,\eta_0}(x)$  is continuous around  $\theta_0$  and  $p_0(x)$  is uniformly bounded away from zero. The second condition is trivially satisfied in the semiparametric logistic regression model by the given form of the density. The boundedness of  $m_{\theta,\eta}$  thus permits the application of lemma 4.2 below which is used to verify condition (52) in the following lemma 4.1. Note that lemma 4.1 and lemma 4.2 are theorem 3.2 and lemma 3.3 in [15], respectively.

**Lemma 4.1**—Assume for any given  $\theta \in \Theta_n$ ,  $\hat{\eta}_\theta$  satisfies  $\mathbb{P}_n m_{\theta,\hat{\eta}_\theta} \geq \mathbb{P}_n m_{\theta,\eta_0}$  for given measurable functions  $x \mapsto m_{\theta,\eta}(x)$ . Assume conditions (51) and (52) below hold for every  $\theta \in \Theta_n$ , every  $\eta \in \mathcal{V}_n$  and every  $\varepsilon > 0$ :

$$P(m_{\theta,\eta} - m_{\theta,\eta_0}) \lesssim -d_\theta^2(\eta, \eta_0) + \|\theta - \theta_0\|^2, \tag{51}$$

$$E^* \sup_{\theta \in \Theta_n, \eta \in \mathcal{V}_n, \|\theta - \theta_0\| < \varepsilon, d_\theta(\eta, \eta_0) < \varepsilon} |\mathbb{G}_n(m_{\theta,\eta} - m_{\theta,\eta_0})| \lesssim \varphi_n(\varepsilon). \tag{52}$$

Suppose that (52) is valid for functions  $\varphi_n$  such that  $\delta \mapsto \varphi_n(\delta)/\delta^\alpha$  is decreasing for some  $\alpha < 2$  and sets  $\Theta_n \times \mathcal{V}_n$  such that  $P(\tilde{\theta} \in \Theta_n, \tilde{\eta}_\theta \in \mathcal{V}_n) \rightarrow 1$ . Then  $d_\theta(\tilde{\eta}_\theta, \eta_0) \leq O_p^*(\delta_n + \|\tilde{\theta} - \theta_0\|)$  for any sequence of positive numbers  $\delta_n$  such that  $\varphi_n(\delta_n) \leq \sqrt{n}\delta_n^2$  for every  $n$ .

Lemma 4.2 below is presented to verify the modulus condition for the continuity of the empirical process in (52). Let  $\mathcal{S}_\delta = \{x \mapsto m_{\theta,\eta}(x) - m_{\theta,\eta_0}(x) : d_\theta(\eta, \eta_0) < \delta, \|\theta - \theta_0\| < \delta\}$  and write

$$K(\delta, \mathcal{S}_\delta, L_2(P)) = \int_0^\delta \sqrt{1 + H_b(\varepsilon, \mathcal{S}_\delta, L_2(P))} d\varepsilon. \tag{53}$$

**Lemma 4.2**—Suppose the functions  $(x, \theta, \eta) \mapsto m_{\theta,\eta}(x)$  are uniformly bounded for  $(\theta, \eta)$  ranging over a neighborhood of  $(\theta_0, \eta_0)$  and that

$$P(m_{\theta,\eta} - m_{\theta_0,\eta_0})^2 \lesssim d_\theta^2(\eta, \eta_0) + \|\theta - \theta_0\|^2.$$

Then condition (52) is satisfied for any functions  $\varphi_n$  such that

$$\varphi_n(\delta) \geq K(\delta, \mathcal{S}_\delta, L_2(P)) \left( 1 + \frac{K(\delta, \mathcal{S}_\delta, L_2(P))}{\delta^2 \sqrt{n}} \right)$$

Consequently, in the conclusion of the above theorem, we may use  $K(\delta, \mathcal{S}_\delta, L_2(P))$  rather than  $\varphi_n(\delta)$ .

We then apply lemma 4.1 to the penalized semiparametric logistic regression model by including  $\lambda$  in  $\theta$ , i.e.  $m_{\theta,\lambda,\eta} = m_{\theta,\eta} - \frac{1}{2}\lambda^2(J^2(\eta) - J^2(\eta_0))$ , in the proof of lemma 4. First, lemma 7.1 in [15] establishes that

$$\|p_{\tilde{\theta}_n, \tilde{\eta}_{\tilde{\theta}_n, \lambda_n}} - p_{\theta_0, \eta_0}\|_2 + \lambda_n J(\tilde{\eta}_{\tilde{\theta}_n, \lambda_n}) = O_p(\lambda_n + \|\tilde{\theta}_n - \theta_0\|) \tag{54}$$

after choosing

$$m_{\theta,\lambda,\eta} = \log \frac{p_{\theta,\eta} + p_{\theta_0,\eta_0}}{2p_{\theta_0,\eta_0}} - \frac{1}{2}\lambda^2(J^2(\eta) - J^2(\eta_0))$$

in lemma 4.1. Note that the map  $\theta \mapsto p_{\theta,\eta_0} f^{W,Z}(w, z)$  is uniformly bounded away from zero at  $\theta = \theta_0$  and continuous around a neighborhood of  $\theta_0$ . Hence  $m_{\theta,\lambda,\eta}$  is well defined. Moreover,  $\mathbb{P}_n m_{\theta,\lambda,\eta} \hat{\eta}_\theta \geq \mathbb{P}_n m_{\theta,\lambda,\eta_0}$  by the inequality that  $((p_{\theta,\eta} + p_{\theta_0,\eta_0})/2p_{\theta_0,\eta_0})^2 \geq (p_{\theta,\eta}/p_{\theta_0,\eta_0})$ . (54) now directly implies (31). For the proof of (30), we need to consider the conclusion of lemma 7.4 (i), which states that

$$\|p_{\theta,\eta} - p_{\theta_0,\eta_0}\|_2 \gtrsim (\|\theta - \theta_0\| \wedge 1 + \|\eta - \eta_0\| \wedge 1) \tag{55}$$

Thus we have proved (30). For (32), we just replace the  $m_{\theta,\lambda,\eta}$  with  $m_{\theta,0,\eta}$  in the proof of lemma 7.1 in [15]. Thus we can show that  $d_\theta(\eta, \eta_0) = \|p_{\theta,\eta} - p_{\theta_0,\eta_0}\|_2$ . By combining lemma 4.2 and (55), we know that  $\|\hat{\eta}_{\tilde{\theta}_n} - \eta_0\|_2 = O_p(\delta_n + \|\tilde{\theta}_n - \theta_0\|)$ , for  $\delta_n$  satisfying

$K(\delta_n, \mathcal{S}_{\delta_n}, L_2(P)) \leq \sqrt{n}\delta_n^2$ . Note that  $K(\delta, \mathcal{S}_\delta, L_2(P))$  is as defined in (53). By similar analysis as used in the proof of lemma 7.1 in [15] and the strengthened assumption on  $\eta$ , we then find that  $K(\delta_n, \mathcal{S}_{\delta_n}, L_2(P)) \lesssim \delta_n^{1-1/2k}$ , which leads to the desired convergence rate given in (32).

**Proof of lemma 5**—The proof of lemma 5 follows that of lemma 2. The smoothness conditions of  $\ell(t, \theta, \eta)$  and its related derivatives can be shown similarly since  $F(\cdot), \tilde{F}(\cdot)$  and  $\dot{F}(\cdot)$  are all uniformly bounded in  $(-\infty, +\infty)$ , and  $h_0(\cdot)$  is intrinsically bounded over  $[0, 1]$ . Note that we can show (12) directly by the following analysis.  $P\ell(\theta_0, \theta_0, \eta)$  can be written as  $P(F(\theta_0 w + \eta_0) - F(\theta_0 w + \eta(z)))(w - h_0(z))$  since  $P\ell_0 = 0$ . Note that  $P(w - h_0(z))\tilde{F}(\theta_0 w + \eta_0(z))(\eta - \eta_0)(z) = 0$ . This implies that  $P\ell(\theta_0, \theta_0, \eta) = P(F(\theta_0 w + \eta_0) - F(\theta_0 w + \eta(z)) + \tilde{F}(\theta_0 w + \eta_0(z))(\eta - \eta_0)(z))(w - h_0(z))$ . However, by the common Taylor expansion, we have  $|F(\theta_0 w + \eta) - F(\theta_0 w + \eta_0) - \tilde{F}(\theta_0 w + \eta_0)(\eta - \eta_0)| \leq \|\tilde{F}\|_\infty |\eta - \eta_0|^2$ . This proves (12).

We next verify the asymptotic equicontinuity conditions, i.e. (15)–(17). For (17), we first apply analysis similar to that used in the proof of lemma 2 to obtain

$$P \left( \frac{\dot{\ell}(\theta_0, \theta_0, \tilde{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}_0}{n^{\frac{1}{4k+2}}(\lambda_n + \|\tilde{\theta}_n - \theta_0\|)} \right)^2 \lesssim O_p \left( n^{-\frac{1}{2k+1}} \right).$$

By lemma 7.1 in [15], we know that  $J(\hat{\eta}_{\tilde{\theta}_n, \lambda_n}) = O_P(1 + \|\tilde{\theta}_n - \theta_0\|/\lambda_n)$  and  $\|\hat{\eta}_{\tilde{\theta}_n, \lambda_n}\|_\infty$  is bounded in probability by a multiple of  $J(\hat{\eta}_{\tilde{\theta}_n, \lambda_n}) + 1$ . Now we construct the set  $\mathfrak{Q}_n$  as follows:

$$\left\{ \frac{\dot{\ell}(\theta_0, \theta_0, \eta) - \dot{\ell}_0}{n^{\frac{1}{4k+2}}(\lambda_n + \|\theta - \theta_0\|)} : J(\eta) \leq C_n(1 + \frac{\|\theta - \theta_0\|}{\lambda_n}), \|\eta\|_\infty \leq C_n(1 + J(\eta)), \|\theta - \theta_0\| < \delta \right\} \cap \left\{ g \in L_2(P) : Pg^2 \leq C_n n^{-\frac{1}{2k+1}} \right\}.$$

Clearly, the probability that the function  $n^{-1/(4k+2)}(\dot{\ell}(\theta_0, \theta_0, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}_0)/(\lambda_n + \|\tilde{\theta}_n - \theta_0\|) \in \mathfrak{Q}_n$  approaches 1 as  $C_n \rightarrow \infty$ . We next show that  $\lim_{n \rightarrow \infty} E^* \|\mathfrak{G}_n\|_{\mathfrak{Q}_n} < \infty$  by T2. Note that  $\dot{\ell}(\theta_0, \theta_0, \eta)$  depends on  $\eta$  in a Lipschitz manner. Consequently, we can bound  $H_B(\varepsilon, \mathfrak{Q}_n, L_2(P))$  by the product of some constant and  $H(\varepsilon, \mathfrak{R}_n, L_2(P))$  in view of T3, where  $\mathfrak{R}_n$  is as defined in the proof of lemma 2. By similar calculations as those performed in lemma 2, we can obtain

$$K(\delta_n, \tilde{\mathfrak{Q}}_n, L_2(P)) \lesssim \lambda_n^{-1/2k} n^{-1/(4k+2)}. \text{ Thus } \lim_{n \rightarrow \infty} E^* \|\mathfrak{G}_n\|_{\mathfrak{Q}_n} < \infty, \text{ and (17) follows.}$$

The proof of (15) and (16) follows arguments quite similar to those used in the proof of lemma 2. In other words, we can show that  $\mathfrak{G}_n(\dot{\ell}(\theta_0, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}_0) = o_P(1 + n^{1/3}\|\tilde{\theta}_n - \theta_0\|) = o_P(1)$  and  $\mathfrak{G}_n(\dot{\ell}_{t,\theta}(\theta_0, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}_{t,\theta}(\theta_0, \theta_0, \eta_0)) = o_P(1 + n^{1/3}\|\tilde{\theta}_n - \theta_0\|)$ .

Next we define  $\mathfrak{V}_n \equiv \{ \ell^{(3)}(\tilde{\theta}_n, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})(x) \}$ . Similar arguments as those used in the proof of lemma 2 can be directly applied to the verification of (18) in this second model. By the form of  $\ell^{(3)}(t, \theta, \eta)$ , the entropy number for  $\mathfrak{V}_n$  is bounded above by that of  $\mathfrak{F}_n \equiv \{ F'(tw + \eta(z)) + (\theta - t)h_0(z) : (t, \theta) \in V_{\theta_0}, \lambda_n J(\eta) \leq C_n, \|\eta\|_\infty \leq C_n(1 + J(\eta)) \}$ . Similarly, we know

$$\sup_Q H(\varepsilon, \overline{\mathfrak{V}}_n, L_1(Q)) \leq \sup_Q H(\varepsilon, \overline{\mathfrak{F}}_n, L_1(Q)) \lesssim ((1 + \lambda_n^{-1})/\varepsilon)^{1/k} = o_p(n). \text{ Moreover, the } \mathfrak{F}_n \text{ are uniformly bounded. This completes the proof for (18). This concludes the proof.}$$

**Proof of lemma 6**—The proof of lemma 6 is analogous to that of lemma 3.

**Lemma 7**—Assuming the assumptions in theorem 1, we have

$$\log pl_{\lambda_n}(\tilde{\theta}_n) = \log pl_{\lambda_n}(\theta_0) + n(\tilde{\theta}_n - \theta_0)^T \mathbb{P}_n \dot{\ell}_0 - \frac{n}{2}(\tilde{\theta}_n - \theta_0)^T \tilde{I}_0(\tilde{\theta}_n - \theta_0) + O_p(g_{\lambda_n}(\|\tilde{\theta}_n - \tilde{\theta}_{\lambda_n}\|)), \tag{56}$$

for any sequence  $\tilde{\theta}_n$  satisfying  $\tilde{\theta}_n = \theta_0 + o_P(1)$ .

**Proof**— $n^{-1}(\log pl_{\lambda_n}(\tilde{\theta}_n) - \log pl_{\lambda_n}(\theta_0))$  is bounded above and below by

$$\mathbb{P}_n(\dot{\ell}(\tilde{\theta}_n, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}(\theta_0, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})) - \lambda_n^2(J^2(\hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - J^2(\eta_{\theta_0}(\tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})))$$

and

$$\mathbb{P}_n(\dot{\ell}(\tilde{\theta}_n, \theta_0, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}(\theta_0, \theta_0, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})) - \lambda_n^2(J^2(\eta_{\tilde{\theta}_n}(\theta_0, \hat{\eta}_{\tilde{\theta}_n, \lambda_n})) - J^2(\hat{\eta}_{\theta_0, \lambda_n})),$$

respectively. By the third order Taylor expansion of  $\tilde{\theta}_n \mapsto \mathbb{P}_n \dot{\ell}(\tilde{\theta}_n, \theta, \eta)$  around  $\theta_0$ , for  $\theta = \tilde{\theta}_n$  and  $\eta = \hat{\eta}_{\tilde{\theta}_n, \lambda_n}$ , (18) and the above empirical no-bias conditions (13) and (14), we can find that the order of the difference between  $\mathbb{P}_n(\dot{\ell}(\tilde{\theta}_n, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}) - \dot{\ell}(\theta_0, \tilde{\theta}_n, \hat{\eta}_{\tilde{\theta}_n, \lambda_n}))$  and  $(\tilde{\theta}_n - \theta_0)^T \mathbb{P}_n \dot{\ell}_0 - (\tilde{\theta}_n - \theta_0)^T (\tilde{I}_0/2)(\tilde{\theta}_n - \theta_0)$  is  $O_P(n^{-1}g_{\lambda_n}(\|\tilde{\theta}_n - \tilde{\theta}_{\lambda_n}\|))$ . Similarly, we have

$$\begin{aligned} \lambda_n^2(J^2(\widehat{\eta}_{\theta_n, \lambda_n}) - J^2(\eta_{\theta_0}(\widetilde{\theta}_n, \widehat{\eta}_{\theta_n, \lambda_n}))) &= -2\lambda_n^2(\widetilde{\theta}_n - \theta_0)^T \int_{\mathcal{Z}} \widehat{\eta}_{\theta_n, \lambda_n}^{(k)} h_0^{(k)} dz + 2\lambda_n^2(\widetilde{\theta}_n - \theta_0)^T \int_{\mathcal{Z}} h_0^{(k)} h_0^{(k)T} dz (\widetilde{\theta}_n - \theta_0) \\ &= O_p(n^{-1} g_{\lambda_n}(\|\widetilde{\theta}_n - \widehat{\theta}_{\lambda_n}\|)) \end{aligned}$$

by Taylor expansion. The last equation holds because of the assumptions (3) and (19). Similar analysis also applies to the lower bound. This proves (56).



**Table 1**Partly Linear Model with  $\lambda_n = n^{-1/3}$  ( $\theta_0 = 1$  and 200 samples)

<b>n</b>	$n^{2/3} PMLE - CM $	$n^{1/6} SE_M - SE_N $	$n^{2/3} L_M - L_N $	$n^{2/3} U_M - U_N $
50	0.8735	1.5007	2.1653	3.4984
100	0.2269	0.9240	0.6927	1.9507
200	0.2565	1.1440	0.7592	0.7182
800	0.0840	0.9539	0.7756	0.5171

**Table 2**Partly Linear Model with  $\lambda_n = n^{-2/5}$  ( $\theta_0 = 1$  and 200 samples)

$n$	$n^{4/5} PMLE - CM $	$n^{3/10} SE_M - SE_N $	$n^{4/5} L_M - L_N $	$n^{4/5} U_M - U_N $
50	0.7866	0.6826	2.3963	2.9725
100	0.8161	0.2389	0.7007	1.0669
200	0.6654	0.5806	0.5614	0.9427
800	0.6032	0.7836	0.1465	0.3782

$n$ , sample size; PMLE, penalized maximum likelihood estimator; CM, empirical mean;  $SE_M$ , estimated standard errors based on MCMC;  $SE_N$ , estimated standard errors based on numerical derivatives;  $L_M$  ( $U_M$ ), lower (upper) bound of the 95% confidence interval based on MCMC;  $L_N$  ( $U_N$ ), lower (upper) bound of the 95% confidence interval based on numerical derivatives.