

Institute of Mathematical Statistics

LECTURE NOTES — MONOGRAPH SERIES

PARTIAL LIKELIHOOD AND ESTIMATING EQUATIONS

Priscilla E. Greenwood
 University of British Columbia
 and
 Wolfgang Wefelmeyer
 University of Siegen

Abstract

Consider a regression model for discrete-time stochastic processes, with a (partially specified) model for the conditional distribution of the response given the covariate and the past observations. Suppose we also have some knowledge about how the parameter of interest affects the conditional distribution of the covariate given the past. We assume that these two model assumptions give rise to two martingale estimating functions, and determine an optimal combination. We indicate for the case of jump processes how our result carries over to continuous time. The resulting estimators are efficient.

1 Introduction

Suppose we know something about how the parameter of interest in a regression model appears both in the conditional distribution of the response given the covariate, and in the distribution of the covariate. How can we exploit this knowledge? Let us illustrate our approach in the case of independent and identically distributed observations (X_i, Y_i) , with X_i the covariate and Y_i the response. In a regression model one usually specifies the conditional distribution of Y given X , either fully, by a parametric model, or partially. An example of a partial specification is a model for the conditional mean of Y given X , say $E(Y|X) = \vartheta X$. More generally, we specify a function $\bar{g}_\vartheta(X, Y)$ such that $E(\bar{g}_\vartheta(X, Y)|X) = 0$. In the example, $\bar{g}_\vartheta(X, Y) = Y - \vartheta X$. We assume a similar partial specification of the distribution of the covariate X , say $Eg_{*\vartheta}(X) = 0$. The two functions \bar{g}_ϑ and $g_{*\vartheta}$ give rise to two estimating equations

$$\sum_{i=1}^n \bar{g}_\vartheta(X_i, Y_i) = 0, \quad \sum_{i=1}^n g_{*\vartheta}(X_i) = 0.$$

¹Work supported by NSERC, Canada.

By the usual Taylor expansion argument, their solutions $\vartheta = \hat{\vartheta}_n$ are asymptotically normal with variances I^{-1} and I_*^{-1} , respectively, where

$$I = \frac{(\mathbb{E}\bar{g}'_{\vartheta})^2}{\mathbb{E}\bar{g}_{\vartheta}^2}, \quad I_* = \frac{(\mathbb{E}g'_{*\vartheta})^2}{\mathbb{E}g_{*\vartheta}^2}.$$

Since $\bar{g}_{\vartheta}(X_i, Y_i)$ and $g_{*\vartheta}(X_i)$ are uncorrelated, the combined estimating equation

$$\sum (w\bar{g}_{\vartheta}(X_i, Y_i) + w_*g_{*\vartheta}(X_i)) = 0$$

leads to an estimator with asymptotic variance

$$\frac{w^2\mathbb{E}\bar{g}_{\vartheta}^2 + w_*^2\mathbb{E}g_{*\vartheta}^2}{(w\mathbb{E}\bar{g}'_{\vartheta} + w_*\mathbb{E}g'_{*\vartheta})^2}.$$

Applying the Schwarz inequality to the denominator, one sees that this variance is minimized for

$$w^{\text{opt}} = \frac{\mathbb{E}\bar{g}'_{\vartheta}}{\mathbb{E}\bar{g}_{\vartheta}^2}, \quad w_*^{\text{opt}} = \frac{\mathbb{E}g'_{*\vartheta}}{\mathbb{E}g_{*\vartheta}^2}.$$

The minimal asymptotic variance is $(I + I_*)^{-1}$. The weights w^{opt} and w_*^{opt} depend on ϑ and, in general, also on other features of the distribution of (X, Y) . In the estimating function, they must be replaced by estimators \hat{w}_n^{opt} and $\hat{w}_{*n}^{\text{opt}}$, say by using empirical estimators for the distributions involved. This does not change the asymptotic variance $(I + I_*)^{-1}$.

Can we do better than using the combined estimating equation? Note that we can multiply $\bar{g}_{\vartheta}(X, Y)$ by a function $\bar{w}(X)$ of X and still have conditional expectation zero,

$$\mathbb{E}(\bar{w}(X)\bar{g}_{\vartheta}(X, Y)|X) = 0.$$

This leads to new estimating equations

$$\sum (\bar{w}(X_i)\bar{g}_{\vartheta}(X_i, Y_i) + w_*g_{*\vartheta}(X_i)) = 0 \tag{1.1}$$

with asymptotic variance $(I_{\bar{w}} + I_*)^{-1}$, where

$$I_{\bar{w}} = \frac{(\mathbb{E}\bar{w}\bar{g}'_{\vartheta})^2}{\mathbb{E}\bar{w}^2\bar{g}_{\vartheta}^2} = \frac{(\mathbb{E}(\bar{w}(X)\mathbb{E}(\bar{g}'_{\vartheta}(X, Y)|X)))^2}{\mathbb{E}(\bar{w}(X)^2\mathbb{E}(\bar{g}_{\vartheta}(X, Y)^2|X))}.$$

Applying again the Schwarz inequality, one sees that $I_{\bar{w}}$ is maximized by

$$\bar{w}^{\text{opt}}(X) = \frac{\mathbb{E}(\bar{g}'_{\vartheta}(X, Y)|X)}{\mathbb{E}(\bar{g}_{\vartheta}(X, Y)^2|X)}.$$

The weight again depends on ϑ and, in general, also on other features of the distribution of (X, Y) and must be estimated, say by using appropriate

nonparametric estimators for the conditional expectations involved. With $\hat{w}_n^{\text{opt}}(X)$ denoting such an estimator, we arrive at the estimating equation

$$\sum \left(\hat{w}_n^{\text{opt}}(X_i) \bar{g}_\vartheta(X_i, Y_i) + \hat{w}_{*n}^{\text{opt}} g_{*\vartheta}(X_i) \right) = 0. \quad (1.2)$$

The asymptotic variance of the estimator corresponding to this equation is $(\bar{I} + I_*)^{-1}$ with

$$\bar{I} = \mathbb{E} \frac{(\mathbb{E}(\bar{g}'_\vartheta(X, Y)|X))^2}{\mathbb{E}(\bar{g}_\vartheta(X, Y)^2|X)}.$$

By the Schwarz inequality (2.10) below, \bar{I} is strictly larger than I unless both conditional expectations do not depend on X .

For the example given above, $\bar{g}_\vartheta(X, Y) = Y - \vartheta X$, we have $\mathbb{E}(\bar{g}'_\vartheta(X, Y)|X) = -X$ and $\mathbb{E}(\bar{g}_\vartheta(X, Y)^2|X) = \mathbb{E}((Y - \vartheta X)^2|X)$, the conditional variance of Y given X .

The estimating equation (1.2) is not only optimal among estimating equations (1.1) but even efficient among *all* (regular) estimators as long as we do not impose additional restrictions on the distribution of (X, Y) which involve ϑ . Let us give a sketch of the argument, referring to Bickel, Klaassen, Ritov and Wellner (1993) for an account of the concepts involved.

The model is described by all distributions $p(dx, dy) = p_*(dx)\bar{p}(x, dy)$ of (X, Y) such that

$$\int \bar{p}(x, dy) \bar{g}_\vartheta(x, y) = \mathbb{E}(\bar{g}_\vartheta(x, Y)|x) = 0 \quad \text{for all } x, \quad (1.3)$$

$$\int p_*(dx) g_{*\vartheta}(x) = \mathbb{E} g_{*\vartheta} = 0 \quad (1.4)$$

if ϑ is true. Introduce a local model by perturbing \bar{p} and p_* as $\bar{p}(x, dy)(1 + n^{-1/2}u\bar{h}(x, y))$ and $p_*(dx)(1 + n^{-1/2}uh_*(x))$ such that the two conditions (1.3) and (1.4) hold (approximately) with ϑ replaced by $\vartheta + n^{-1/2}u$:

$$\int \bar{p}(x, dy) \left(1 + n^{-1/2}u\bar{h}(x, y) \right) \bar{g}_{\vartheta+n^{-1/2}u}(x, y) = 0 \quad \text{for all } x,$$

$$\int p_*(dx) \left(1 + n^{-1/2}uh_*(x) \right) g_{*,\vartheta+n^{-1/2}u}(x) = 0.$$

Then \bar{h} and h_* must fulfill

$$\int \bar{p}(x, dy) \left(\bar{h}(x, y) \bar{g}_\vartheta(x, y) + \bar{g}'_\vartheta(x, y) \right) = 0 \quad \text{for all } x, \quad (1.5)$$

$$\int p_*(dx) \left(h_*(x) g_{*\vartheta}(x) + g'_{*\vartheta}(x) \right) = 0. \quad (1.6)$$

The perturbed p is approximately

$$p(dx, dy) \left(1 + n^{-1/2}u \left(\bar{h}(x, y) + h_*(x) \right) \right).$$

This means that the *tangent space* in the sense of Bickel et al. (1993, p. 50, Definition 2) consists of the functions $u(\bar{h} + h_*)$ with \bar{h} and h_* fulfilling (1.5) and (1.6). We view the parameter ϑ as a function of p and determine its *canonical gradient* in the sense of Bickel et al. (1993, p. 58). This is a function $\dot{\nu}$ in the tangent space such that

$$n^{1/2}(\vartheta + n^{-1/2}u - \vartheta) = u = uE\dot{\nu}(\bar{h} + h_*) \quad \text{for all } \bar{h}, h_* \text{ fulfilling (1.5), (1.6).}$$

According to Bickel et al. (1993, p. 63, Theorem 2B, and p. 65, Theorem 1A), an estimator $\hat{\vartheta}_n$ is regular and efficient if and only if

$$n^{1/2}(\hat{\vartheta}_n - \vartheta) = n^{-1/2} \sum_{i=1}^n \dot{\nu}(X_i, Y_i) + o_P(1). \quad (1.7)$$

In particular, a lower bound for the asymptotic variance of regular estimators of ϑ is $E\dot{\nu}^2$.

Since the tangent space is generated by the affine space of functions $\bar{h} + h_*$ with \bar{h} and h_* fulfilling (1.5) and (1.6), we can write the canonical gradient as $\dot{\nu} = (E(\bar{s} + s_*)^2)^{-1}(\bar{s} + s_*)$, where $\bar{s} + s_*$ is the *optimal score function*, minimizing $E(\bar{h} + h_*)^2$ over all \bar{h} and h_* fulfilling (1.5) and (1.6). In particular, the lower bound for the asymptotic variance of regular estimators can be written $1/E(\bar{s} + s_*)^2$. The function $\bar{s} + s_*$ is characterized by

$$E(\bar{s} + s_*)(\bar{h} + h_*) = E(\bar{s} + s_*)^2 \quad \text{for all } \bar{h}, h_* \text{ fulfilling (1.5), (1.6).}$$

Since \bar{h} and h_* are orthogonal, this is equivalent to

$$E\bar{s}\bar{h} = E\bar{s}^2 \quad \text{for all } \bar{h} \text{ fulfilling (1.5),} \quad Es_*h_* = Es_*^2 \quad \text{for all } h_* \text{ fulfilling (1.6).}$$

One easily checks that the solution is

$$\bar{s}(x, y) = -\bar{w}^{\text{opt}}(x)\bar{g}_\vartheta(x, y), \quad s_*(x) = -w_*^{\text{opt}}g_{*\vartheta}(x).$$

By the usual Taylor series argument, the solution $\vartheta = \hat{\vartheta}_n$ of the optimal estimating equation (1.2) is seen to fulfill

$$n^{1/2}(\hat{\vartheta}_n - \vartheta) = \left(E(\bar{s} + s_*)^2\right)^{-1} n^{-1/2} \sum_{i=1}^n (\bar{s}(X_i, Y_i) + s_*(X_i)) + o_P(1).$$

By the characterization (1.7), this means that this estimator is efficient.

In Sections 2 and 3 we show how the calculation of the optimal estimating function carries over to ergodic discrete-time stochastic processes and jump processes, respectively. Efficiency also carries over, but we will not give the details. All results extend immediately to vectors ϑ and vector-valued functions \bar{g}_ϑ and $g_{*\vartheta}$. We do not give precise regularity conditions for our results.

2 Discrete-time stochastic processes

Suppose we observe a stochastic process (X_i, Y_i) at times $i = 1, \dots, n$. The law of the process is determined by the conditional distributions $p_i(dx, dy)$ of (X_i, Y_i) given the past observations. Here and in the following, we suppress the dependence of p_i and similar objects on the past, $(X_1, Y_1), \dots, (X_{i-1}, Y_{i-1})$. As in the i.i.d. case considered in Section 1, we describe a regression model by (partial) specifications of (1) the conditional distribution of the response given the present value of the covariate and now also the past observations, and of (2) the conditional distribution of the covariate, now also given the past. We factor p_i into marginal and conditional,

$$p_i(dx, dy) = p_{*i}(dx)\bar{p}_i(x, dy), \quad (2.1)$$

and specify two functions $\bar{g}_{i\vartheta}(x, y)$ and $g_{*i\vartheta}(x)$, possibly depending on the past, such that

$$\begin{aligned} \bar{E}^x \bar{g}_{i\vartheta} &= \int \bar{p}_i(x, dy) \bar{g}_{i\vartheta}(x, y) = 0 \text{ for all } x, \\ E_{*i} g_{*i\vartheta} &= \int p_{*i}(dx) g_{*i\vartheta}(x) = 0. \end{aligned}$$

They give rise to estimating equations

$$\sum_{i=1}^n \bar{g}_{i\vartheta}(X_i, Y_i) = 0, \quad \sum_{i=1}^n g_{*i\vartheta}(X_i) = 0.$$

How can we combine them in an optimal way? Our result holds for geometrically ergodic processes and under appropriate smoothness and moment conditions which can be seen from the sketch of the proof.

Result 1. *From estimating equations of the form*

$$\sum_{i=1}^n (\bar{w}_i(X_i) \bar{g}_{i\vartheta}(X_i, Y_i) + w_{*i} g_{*i\vartheta}(X_i)) = 0, \quad (2.2)$$

an estimator with minimal asymptotic variance is obtained using weights which are consistent estimators of

$$\bar{w}_i(X_i) = \bar{E}_i^{X_i} \bar{g}'_{i\vartheta} / \bar{E}_i^{X_i} \bar{g}_{i\vartheta}^2, \quad (2.3)$$

$$w_{*i} = E_{*i} g'_{*i\vartheta} / E_{*i} g_{*i\vartheta}^2. \quad (2.4)$$

The estimator is asymptotically normal. Its asymptotic variance is the limit of

$$n \left(\sum_{i=1}^n E_{*i} \frac{(\bar{E}_i^{X_i} \bar{g}'_{i\vartheta})^2}{\bar{E}_i^{X_i} \bar{g}_{i\vartheta}^2} + \sum_{i=1}^n \frac{(E_{*i}^{X_i} g'_{*i\vartheta})^2}{E_{*i}^{X_i} g_{*i\vartheta}^2} \right)^{-1}. \quad (2.5)$$

Sketch of proof. To simplify the notation, we introduce $w = (\bar{w}, w_*)$ and $g_\vartheta = (\bar{g}_\vartheta, g_{*\vartheta})$, and write the estimating equation (2.2) as

$$\sum w_i g_{i\vartheta} = \sum \bar{w}_i \bar{g}_{i\vartheta} + \sum w_{*i} g_{*\vartheta} = 0.$$

Let $\hat{\vartheta}$ be a solution. Under appropriate differentiability conditions, a Taylor expansion gives

$$0 = \sum w_i g_{i\hat{\vartheta}} \approx \sum w_i g_{i\vartheta} + (\hat{\vartheta} - \vartheta) \sum w_i g'_{i\vartheta}. \quad (2.6)$$

Then

$$n^{1/2}(\hat{\vartheta} - \vartheta) \approx -n^{-1/2} \sum w_i g_{i\vartheta} / \frac{1}{n} \sum w_i g'_{i\vartheta}. \quad (2.7)$$

Conditionally on the past, the martingale increments $\bar{w}_i(X_i) \bar{g}_{i\vartheta}(X_i, Y_i)$ and $w_{*i} g_{*\vartheta}(X_i)$ are orthogonal:

$$\mathbb{E}_i \bar{w}_i \bar{g}_{i\vartheta} w_{*i} g_{*\vartheta} = w_{*i} \int p_{*i}(dx) g_{*\vartheta}(x) \bar{w}_i(x) \int \bar{p}_i(x, dy) \bar{g}_{i\vartheta}(x, y) = 0. \quad (2.8)$$

Introduce an inner product

$$(v, w) = \sum \mathbb{E}_{*i} \bar{v}_i \bar{w}_i + \sum v_{*i} w_{*i}$$

with corresponding norm $\|w\|^2 = (w, w)$. Interpret products vw of vectors componentwise,

$$vw = (\bar{v}_1 \bar{w}_1, \dots, \bar{v}_n \bar{w}_n, v_{*1} w_{*1}, \dots, v_{*n} w_{*n}).$$

Consider first the numerator in (2.7). With (2.8), the predictable quadratic variation of $\sum w_i g_{i\vartheta}$ is $(w^2, v) = \|wv^{1/2}\|^2$, where

$$\bar{v}_i(x) = \bar{E}_i^x \bar{g}_{i\vartheta}^2, \quad v_{*i} = \mathbb{E}_{*i} g_{*\vartheta}^2.$$

Consider now the denominator in (2.7). The compensator of $\sum w_i g'_{i\vartheta}$ is (w, m) , where

$$\bar{m}_i(x) = \bar{E}_i^x \bar{g}'_{i\vartheta}, \quad m_{*i} = \mathbb{E}_{*i} g'_{*\vartheta}.$$

Since $\sum w_i g'_{i\vartheta} - (w, m)$ is a martingale, $\frac{1}{n} (\sum w_i g'_{i\vartheta} - (w, m))$ is asymptotically negligible, and we may replace $\frac{1}{n} \sum w_i g'_{i\vartheta}$ by $\frac{1}{n} (w, m)$. If the process is ergodic, $\frac{1}{n} (w, m)$ is asymptotically constant. Hence the predictable quadratic variation of $n^{1/2}(\hat{\vartheta} - \vartheta)$ is approximately $n \|wv^{1/2}\|^2 / (w, m)^2$. By the Schwarz inequality,

$$|(w, m)| = |(wv^{1/2}, mv^{-1/2})| \leq \|wv^{1/2}\| \|mv^{-1/2}\|. \quad (2.9)$$

In other words,

$$\frac{\|wv^{1/2}\|^2}{(w, m)^2} \geq \frac{1}{\|mv^{-1/2}\|^2} = \frac{\|mv^{-1} \cdot v^{1/2}\|^2}{(mv^{-1}, m)^2}.$$

Hence $\|wv^{1/2}\|^2/(w, m)^2$ is minimized by $w = mv^{-1}$, and the minimum is

$$\|mv^{-1/2}\|^{-2} = \left(\sum E_{*i} \frac{\bar{m}_i^2}{\bar{v}_i} + \sum \frac{m_{*i}^2}{v_{*i}} \right)^{-1}.$$

By an appropriate martingale central limit theorem, $n^{1/2}(\hat{\vartheta} - \vartheta)$ is asymptotically normal with variance equal to the limit of $n\|mv^{-1/2}\|^{-2}$, and the assertion follows. \square

Efficiency of the estimator based on the optimal estimating equation can be proved by an approach similar to that outlined in Section 1 for the i.i.d. case.

If we use *predictors* to estimate $\bar{w}_i(x)$ and w_{*i} , i.e. estimators involving only the past observations $(X_1, Y_1), \dots, (X_{i-1}, Y_{i-1})$, then the optimal estimating function is a martingale.

We may allow weights \bar{w}_i, w_{*i} to depend on ϑ . Then the derivative of $\sum w_i g_{i\vartheta}$ in the expansion (2.6) has a second term $\sum w'_i g_{i\vartheta}$. It is asymptotically negligible since the $g_{i\vartheta}$ are martingale increments.

Remark 1. Usually one takes *predictable* weights to combine two martingale estimating functions; e.g. Heyde (1987). For the estimating functions $\sum \bar{g}_{i\vartheta}(X_i, Y_i)$ and $\sum g_{*i\vartheta}(X_i)$ this would mean using weights \bar{w}_i rather than $\bar{w}_i(X_i)$. Then the best weights would be

$$\bar{w}_i = E_i \bar{g}'_{i\vartheta} / E_i \bar{g}_{i\vartheta}^2$$

and w_{*i} as above, and the minimal asymptotic variance would be the limit of an expression of the form (2.5) with $E_{*i} \left((\bar{E}_i^{X_i} \bar{g}'_{i\vartheta})^2 / \bar{E}_i^{X_i} \bar{g}_{i\vartheta}^2 \right)$ replaced by the simpler $(E_i \bar{g}'_{i\vartheta})^2 / E_i \bar{g}_{i\vartheta}^2$. The resulting variance is, in general, larger than our minimum variance (2.5) because

$$\frac{(E f)^2}{E g^2} \leq E \frac{(E^X f)^2}{E^X g^2}. \tag{2.10}$$

This inequality follows from the Schwarz inequality:

$$\begin{aligned} (E f)^2 &= (E E^X f)^2 = \left(E \frac{E^X f}{(E^X g^2)^{1/2}} (E^X g^2)^{1/2} \right)^2 \\ &\leq E \frac{(E^X f)^2}{E^X g^2} \cdot E E^X g^2 = E \frac{(E^X f)^2}{E^X g^2} \cdot E g^2. \end{aligned}$$

Remark 2. The weight \bar{w}_i depends only on \bar{p}_i and $\bar{g}_{i\vartheta}$, and w_{*i} depends only on p_{*i} and $g_{*i\vartheta}$. This is due to orthogonality (2.8). Indeed, the weights (2.3) are optimal for estimating functions of the form $\sum \bar{w}_i(X_i)\bar{g}_{i\vartheta}(X_i, Y_i)$, and the weights (2.4) are optimal for estimating functions of the form $\sum w_{*i}g_{*i\vartheta}(X_i)$.

Remark 3. Suppose we have a parametric model $\bar{p}_{i\vartheta}(X_i, dy)$ for the conditional distribution of the response Y_i given the present covariate X_i and the past observations. Differentiating under the integral, we obtain

$$0 = \left(\bar{E}_{i\vartheta}^x \bar{g}_{i\vartheta} \right)' = \bar{E}_{i\vartheta}^x \bar{g}'_{i\vartheta} + \bar{E}_{i\vartheta}^x \bar{\ell}'_{i\vartheta} \bar{g}_{i\vartheta},$$

where

$$\bar{\ell}'_{i\vartheta}(x, y) = \partial_{\tau=\vartheta} \frac{d\bar{p}_{i\tau}(x, \cdot)}{d\bar{p}_{i\vartheta}(x, \cdot)}(y).$$

Hence, by the Schwarz inequality, $\left(\bar{E}_{i\vartheta}^x \bar{g}'_{i\vartheta} \right)^2 / \bar{E}_{i\vartheta}^x \bar{g}_{i\vartheta}^2$ is maximal for $\bar{g}_{i\vartheta} = \bar{\ell}'_{i\vartheta}$, and the optimal weight (2.3) for $\bar{g}_{i\vartheta} = \bar{\ell}'_{i\vartheta}$ is $\bar{w}_i(X_i) = -1$.

In particular, the estimating function $\sum \bar{\ell}'_{i\vartheta}(X_i, Y_i)$ is optimal among estimating functions $\sum \bar{w}_i(X_i)\bar{g}_{i\vartheta}(X_i, Y_i)$. The optimal estimating function is the *partial score function*, i.e., the derivative of the *partial likelihood ratio* of Cox (1975),

$$\prod \frac{d\bar{p}_{i\tau}(X_i, \cdot)}{d\bar{p}_{i\vartheta}(X_i, \cdot)}(Y_i).$$

Hence the optimal estimating function gives the *maximum partial likelihood estimator*. If the observations (X_i, Y_i) are independent, the partial likelihood ratio is the *conditional likelihood ratio* for Y_1, \dots, Y_n given the covariates X_1, \dots, X_n .

Similarly, if there is a parametric model $p_{*i\vartheta}$ for the distribution of the covariate X_i given the past, the optimal $g_{*i\vartheta}$ is

$$g_{*i\vartheta}(x) = \partial_{\tau=\vartheta} \frac{dp_{*i\tau}(x)}{dp_{*i\vartheta}(x)} =: \ell'_{*i\vartheta}(x),$$

and the optimal weight (2.4) is -1 .

Moreover, if there is a *fully* specified parametric model

$$p_{i\vartheta}(dx, dy) = p_{*i\vartheta}(dx)\bar{p}_{i\vartheta}(x, dy),$$

then the likelihood ratio can be written

$$\prod \frac{dp_{*i\tau}(X_i)}{dp_{*i\vartheta}(X_i)} \prod \frac{d\bar{p}_{i\tau}(X_i, \cdot)}{d\bar{p}_{i\vartheta}(X_i, \cdot)}(Y_i), \quad (2.11)$$

and the optimal estimating function is the *score function* $\sum \ell'_{i\vartheta}(X_i, Y_i)$ with

$$\ell'_{i\vartheta}(x, y) := \partial_{\tau=\vartheta} \frac{dp_{i\tau}}{dp_{i\vartheta}}(x, y) = \bar{\ell}'_{i\vartheta}(x, y) + \ell'_{*i\vartheta}(x).$$

Hence the optimal estimator is the maximum likelihood estimator, and its asymptotic variance is the limit of

$$n \left(\sum E_i \bar{\ell}'_{i\vartheta}{}^2 + \sum E_{*i} \ell'_{*i\vartheta}{}^2 \right)^{-1}.$$

Remark 4. For discrete-time processes, it is common to model the conditional distribution of the response given the past *and* the present value of the covariate. In the continuous-time setting of Section 3 one usually models the conditional distribution of the response given only the past. This is just a convention: We may consider X_{i-1} rather than X_i the ‘present’ covariate of the response.

3 Jump processes

Suppose we observe a jump process $(X, Y) = (X_s, Y_s)_{s \geq 0}$ on a finite time interval $[0, t]$. The corresponding multivariate point process is given by the jump measure

$$\mu(ds, dx, dy) = \sum_{s: (\Delta X_s, \Delta Y_s) \neq 0} \varepsilon_{(s, \Delta X_s, \Delta Y_s)}(ds, dx, dy).$$

The law of the process is determined by the compensator of the jump measure. Assume, for simplicity, that the compensator has the form $K_s(dx, dy)ds$, so that there are no time points at which the process has a positive probability of jumping. We can write $K_s(dx, dy) = a_s p_s(dx, dy)$ with p_s a *probability* measure, the jump size distribution at time s given the past, and a_s the jump intensity. For the theory of continuous-time processes and limit theorems we refer to Jacod and Shiryaev (1987).

The multivariate point process corresponding to the response process Y is

$$\mu^Y(ds, dy) = \sum_{s: \Delta Y_s \neq 0} \varepsilon_{(s, \Delta Y_s)}(ds, dy).$$

A regression model is given by a (partial) specification of the compensator of μ^Y , say $\bar{K}_s(dy)ds$. As noted in Remark 4, this is not exactly analogous to the discrete-time case. We specify a predictable function $\bar{g}_{s\vartheta}(y)$ such that

$$\bar{K}_s \bar{g}_{s\vartheta} = \int \bar{K}_s(dy) \bar{g}_{s\vartheta}(y) = 0$$

and obtain a martingale estimating equation

$$\sum_{s \leq t: \Delta Y_s \neq 0} \bar{g}_{s\vartheta}(\Delta Y_s) = 0.$$

We want to assume a similar partial specification of the distribution of the covariate process X . It will be based on a factorization of $\bar{K}_s(dx, dy)$ analogous to the factorization (2.1) of the distribution $p_i(dx, dy)$. We must take into account the possibility that X jumps while Y does not. Following Arjas and Haara (1984) and Greenwood and Wefelmeyer (1996), we write

$$K_s(dx, dy) = K_{-0s}(dx, dy) + K_{*s}(dx)\varepsilon_0(dy),$$

where $K_{-0s}(dx, dy)$ does not charge the subspace described by $y = 0$. Then K_{*s} governs those jumps of X that do not occur simultaneously with jumps of Y . As in (2.1), but with the roles of X and Y interchanged, we factor

$$K_{-0s}(dx, dy) = \bar{K}_s(dy)K_{-*s}(dx, y). \quad (3.1)$$

We note that $K_{-*s}(dx, y)$ is a *probability* measure, the conditional distribution of the jump size of the covariate given a jump of size y of the response, and given the past. Additional specifications of the model may now be given by predictable functions $g_{*s\vartheta}(x)$ and $g_{-*s\vartheta}(x, y)$ such that

$$\begin{aligned} K_{*s}g_{*s\vartheta} &= \int K_{*s}(dx)g_{*s\vartheta}(x) = 0, \\ K_{-*s}^y g_{-*s\vartheta} &= \int K_{-*s}(dx, y)g_{-*s\vartheta}(x, y) = 0 \quad \text{for all } y. \end{aligned}$$

They give rise to additional martingale estimating equations

$$\begin{aligned} \sum_{s \leq t: \Delta X_s \neq 0, \Delta Y_s = 0} g_{*s\vartheta}(\Delta X_s) &= 0, \\ \sum_{s \leq t: \Delta Y_s \neq 0} g_{-*s\vartheta}(\Delta X_s, \Delta Y_s) &= 0. \end{aligned}$$

How can we combine the three estimating functions in an optimal way? Again, our result holds for geometrically ergodic processes under appropriate smoothness and moment conditions which can be seen from the sketch of the proof.

Result 2. *From estimating equations of the form*

$$\begin{aligned} &\sum_{s \leq t: \Delta Y_s \neq 0} \bar{w}_s \bar{g}_{s\vartheta}(\Delta Y_s) \\ &+ \sum_{s \leq t: \Delta X_s \neq 0, \Delta Y_s = 0} w_{*s} g_{*s\vartheta}(\Delta X_s) \\ &+ \sum_{s \leq t: \Delta Y_s \neq 0} w_{-*s}(\Delta Y_s) g_{-*s\vartheta}(\Delta X_s, \Delta Y_s) = 0, \end{aligned} \quad (3.2)$$

an estimator with minimal asymptotic variance is obtained using weights which are consistent estimators of

$$\bar{w}_s = \bar{K}_s \bar{g}'_{s\vartheta} / \bar{K}_s \bar{g}_{s\vartheta}^2, \quad (3.3)$$

$$w_{*s} = K_{*s} g'_{*s\vartheta} / K_{*s} g_{*s\vartheta}^2, \quad (3.4)$$

$$w_{-*s}(y) = K_{-*s}^y g'_{-*s\vartheta} / K_{-*s}^y g_{-*s\vartheta}^2. \quad (3.5)$$

The estimator is asymptotically normal. Its asymptotic variance is the limit of

$$t \left(\int \frac{(\bar{K}_s \bar{g}'_{s\vartheta})^2}{\bar{K}_s \bar{g}_{s\vartheta}^2} ds + \int \frac{(K_{*s} g'_{*s\vartheta})^2}{K_{*s} g_{*s\vartheta}^2} ds + \int \bar{K}_s(dy) \frac{(K_{-*s}^y g'_{-*s\vartheta})^2}{K_{-*s}^y g_{-*s\vartheta}^2} ds \right)^{-1}. \quad (3.6)$$

Sketch of proof. To simplify the notation, we introduce $w = (\bar{w}, w_*, w_{-*})$ and $g_\vartheta = (\bar{g}_\vartheta, g_{*\vartheta}, g_{-*\vartheta})$, and write the estimating equation (3.2) as

$$\sum w_s g_{s\vartheta} = \sum \bar{w}_s \bar{g}_{s\vartheta} + \sum w_{*s} g_{*s\vartheta} + \sum w_{-*s} g_{-*s\vartheta} = 0.$$

Let $\hat{\vartheta}$ be a solution. Under appropriate differentiability conditions, a Taylor expansion gives

$$0 = \sum w_s g_{s\hat{\vartheta}} \approx \sum w_s g_{s\vartheta} + (\hat{\vartheta} - \vartheta) \sum w_s g'_{s\vartheta}. \quad (3.7)$$

Then

$$t^{1/2}(\hat{\vartheta} - \vartheta) \approx -t^{-1/2} \sum w_s g_{s\vartheta} / \frac{1}{t} \sum w_s g'_{s\vartheta}. \quad (3.8)$$

The martingale $\sum w_{*s} g_{*s\vartheta}$ is orthogonal to the martingales $\sum \bar{w}_s \bar{g}_{s\vartheta}$ and $\sum w_{-*s} g_{-*s\vartheta}$ because it lives on time points s with $\Delta Y_s = 0$ while the two other martingales do not jump at these time points. Because $K_{-0s}(dx, dy)$ does not charge $y = 0$, we may and will assume that $g_{-*s\vartheta}(x, 0) = 0$. Then

$$\begin{aligned} K_s \bar{w}_s \bar{g}_{s\vartheta} w_{-*s} g_{-*s\vartheta} &= K_{-0s} \bar{w}_s \bar{g}_{s\vartheta} w_{-*s} g_{-*s\vartheta} \\ &= \bar{w}_s \int \bar{K}_s(dy) \bar{g}_{s\vartheta}(y) w_{-*s}(y) \int K_{-*s}(dx, y) g_{-*s\vartheta}(x, y) = 0. \end{aligned} \quad (3.9)$$

Hence $\sum \bar{w}_s \bar{g}_{s\vartheta}$ and $\sum w_{-*s} g_{-*s\vartheta}$ are also orthogonal. Introduce an inner product

$$(v, w) = \int_0^t \bar{v}_s \bar{w}_s + \int_0^t v_{*s} w_{*s} + \int_0^t \int \bar{K}_s(dy) v_{-*s}(y) w_{-*s}(y) ds$$

with corresponding norm $\|w\|^2 = (w, w)$. Consider first the numerator in (3.8). With (3.9) and orthogonality of $\sum w_{*s} g_{*s\vartheta}$ and $\sum \bar{w}_s \bar{g}_{s\vartheta}$, the predictable quadratic variation of $\sum w_s g_{s\vartheta}$ is $(w^2, v) = \|w v^{1/2}\|^2$, where

$$\bar{v}_s = \bar{K}_s \bar{g}_{s\vartheta}^2, \quad v_{*s} = K_{*s} g_{*s\vartheta}^2, \quad v_{-*s}(y) = K_{-*s}^y g_{-*s\vartheta}^2.$$

Consider now the denominator in (3.8). The compensator of $\sum w_s g_{s\vartheta}$ is (w, m) , where

$$\bar{m}_s = \bar{K}_s \bar{g}'_{s\vartheta}, \quad m_{*s} = K_{*s} g'_{*s\vartheta}, \quad m_{-*s}(y) = K_{-*s}^y g'_{-*s\vartheta}.$$

Since $\sum w_s g'_{s\vartheta} - (w, m)$ is a martingale, $\frac{1}{t} (\sum w_s g'_{s\vartheta} - (w, m))$ is asymptotically negligible, and we may replace $\frac{1}{t} \sum w_s g'_{s\vartheta}$ by $\frac{1}{t} (w, m)$. If the process is ergodic, $\frac{1}{t} (w, m)$ is asymptotically constant. Hence the predictable quadratic variation of $t^{1/2}(\hat{\vartheta} - \vartheta)$ is approximately $t \|wv^{1/2}\|^2 / (w, m)^2$. By the inequality (2.9), this is minimized by $w = mv^{-1}$, and the minimum is

$$t \|mv^{-1/2}\|^{-2} = t \left(\int_0^t \frac{\bar{m}_s^2}{\bar{v}_s} ds + \int_0^t \frac{m_{*s}^2}{v_{*s}} ds + \int_0^t \int \bar{K}_s(dy) \frac{m_{-*s}(y)^2}{v_{-*s}(y)} ds \right)^{-1}.$$

By an appropriate martingale central limit theorem, $t^{1/2}(\hat{\vartheta} - \vartheta)$ is asymptotically normal with variance equal to the limit of $t \|mv^{-1/2}\|^{-2}$, and the assertion follows. \square

As in Remark 2, the weight (3.3) is optimal for estimating equations

$$\sum_{s \leq t: \Delta Y_s \neq 0} \bar{w}_s \bar{g}_{s\vartheta}(\Delta Y_s) = 0,$$

and the weights (3.4) and (3.5) have analogous optimality properties on their own.

Remark 5. Suppose we have a parametric model $\bar{K}_{s\vartheta}(dy)$ for the compensator of the jump measure μ^Y of the response. Write

$$\bar{V}_{s\vartheta\tau} = \frac{d\bar{K}_{s\tau}}{d\bar{K}_{s\vartheta}}, \quad \bar{\ell}'_{s\vartheta} = \partial_{\tau=\vartheta} \bar{V}_{s\vartheta\tau}.$$

When the intensity $\bar{K}_{s\vartheta}(\mathbf{R})$ of the response depends on ϑ , then $\bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta}$ will not be zero in general. This differs from the discrete-time case. Differentiating under the integral, we obtain

$$0 = (\bar{K}_{s\vartheta} \bar{g}_{s\vartheta})' = \bar{K}_{s\vartheta} \bar{g}'_{s\vartheta} + \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta} \bar{g}_{s\vartheta} = \bar{K}_{s\vartheta} \bar{g}'_{s\vartheta} + \bar{K}_{s\vartheta} (\bar{\ell}'_{s\vartheta} - \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta}) \bar{g}_{s\vartheta}.$$

Using the Schwarz inequality, we see that

$$\frac{(\bar{K}_{s\vartheta} \bar{g}'_{s\vartheta})^2}{\bar{K}_{s\vartheta} \bar{g}_{s\vartheta}^2} = \frac{(\bar{K}_{s\vartheta} (\bar{\ell}'_{s\vartheta} - \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta}) \bar{g}_{s\vartheta})^2}{\bar{K}_{s\vartheta} \bar{g}_{s\vartheta}^2}$$

is maximal for $\bar{g}_{s\vartheta} = \bar{\ell}'_{s\vartheta} - \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta}$, and then the optimal weight is $\bar{w}_s = -1$.

In particular,

$$\sum_{s \leq t: \Delta Y_s \neq 0} \left(\bar{\ell}'_{s\vartheta}(\Delta Y_s) - \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta} \right) \quad (3.10)$$

is optimal among estimating functions of the form

$$\sum_{s \leq t: \Delta Y_s \neq 0} \bar{w}_s \bar{g}_{s\vartheta}(\Delta Y_s),$$

and the asymptotic variance of the optimal estimator is the limit of

$$t \left(\int_0^t \bar{K}_{s\vartheta} (\bar{\ell}'_{s\vartheta} - \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta})^2 ds \right)^{-1}.$$

As in Remark 3, the optimal estimating function turns out to be the partial score function in the following sense. A *partial likelihood ratio* for jump processes was introduced by Arjas and Haara (1984) as

$$\prod_{s \leq t: \Delta Y_s \neq 0} \bar{V}_{s\vartheta\tau}(\Delta Y_s) \exp \left(- \int_0^t \bar{K}_{s\vartheta} (\bar{V}_{s\vartheta\tau} - 1) ds \right).$$

See also Andersen, Borgan, Gill and Keiding (1993) and, for general semimartingales, Jacod (1987) and (1990). The *partial score function* is the derivative of the partial likelihood ratio at $\tau = \vartheta$. Using $\bar{V}_{s\vartheta\vartheta} = 1$ we see that the derivative equals (3.10).

Remark 6. Suppose we have a parametric model $K_{*s\vartheta}(dx)ds$ for the compensator of the jump measure of those jumps of the covariate X that do not occur simultaneously with jumps of the response Y . Write

$$V_{*s\vartheta\tau} = \frac{dK_{*s\tau}}{dK_{*s\vartheta}}, \quad \ell'_{*s\vartheta} = \partial_{\tau=\vartheta} V_{*s\vartheta\tau}.$$

As in Remark 5, the best $g_{*s\vartheta}$ is $\ell'_{*s\vartheta} - K_{*s\vartheta} \ell'_{*s\vartheta}$, and then the optimal weight (3.4) is $w_{*s} = -1$.

In particular,

$$\sum_{s \leq t: \Delta X_s \neq 0, \Delta Y_s = 0} \left(\ell'_{*s\vartheta}(\Delta X_s) - K_{*s\vartheta} \ell'_{*s\vartheta} \right) \quad (3.11)$$

is optimal among estimating functions of the form

$$\sum_{s \leq t: \Delta X_s \neq 0, \Delta Y_s = 0} w_{*s} g_{*s\vartheta}(\Delta X_s).$$

Remark 7. Suppose we have a parametric model for the compensator $K_{-*s\vartheta}(dx, y)$ of the conditional jump size distribution at time s of X given a jump $Y = y$ and the past. Write

$$V_{-*s\vartheta\tau}(x, y) = \frac{dK_{-*s\tau}(\cdot, y)}{dK_{-*s\vartheta}(\cdot, y)}, \quad \ell'_{-*s\vartheta} = \partial_{\tau=\vartheta} V_{-*s\vartheta\tau}.$$

Since $K_{-*s\vartheta}(dx, y)$ is a *probability* measure, we have $K_{-*s\vartheta}^y \ell'_{-*s\vartheta} = 0$. As in Remark 5, the best $g_{-*s\vartheta}$ is $\ell'_{-*s\vartheta}$, and then the optimal weight (3.5) is $w_{-*s} = -1$.

In particular,

$$\sum_{s \leq t: \Delta Y_s \neq 0} \ell'_{-*s\vartheta}(\Delta X_s, \Delta Y_s) \quad (3.12)$$

is optimal among estimating functions of the form

$$\sum_{s \leq t: \Delta Y_s \neq 0} w_{-*s} g_{-*s\vartheta}(\Delta X_s, \Delta Y_s).$$

Remark 8. Suppose we have a fully parametric model $\bar{K}_{s\vartheta}$, $K_{*s\vartheta}$, $K_{-*s\vartheta}$. According to Result 2 and Remarks 5 to 7, the best estimating equation is

$$\begin{aligned} & \sum_{s \leq t: \Delta Y_s \neq 0} \left(\bar{\ell}'_{s\vartheta}(\Delta Y_s) - \bar{K}_{s\vartheta} \bar{\ell}'_{s\vartheta} \right) \\ & + \sum_{s \leq t: \Delta X_s \neq 0, \Delta Y_s = 0} \left(\ell'_{*s\vartheta}(\Delta X_s) - K_{*s\vartheta} \ell'_{*s\vartheta} \right) \\ & + \sum_{s \leq t: \Delta Y_s \neq 0} \ell'_{-*s\vartheta}(\Delta X_s, \Delta Y_s) = 0. \end{aligned}$$

To show that this gives the maximum likelihood estimator, we recall a representation of Greenwood and Wefelmeyer (1996) of the likelihood ratio,

$$\begin{aligned} & \prod_{s \leq t: \Delta Y_s \neq 0} \bar{V}_{s\vartheta\tau}(\Delta Y_s) \exp \left(- \int_0^t \bar{K}_{s\vartheta} (\bar{V}_{s\vartheta\tau} - 1) ds \right) \quad (3.13) \\ & \prod_{s \leq t: \Delta X_s \neq 0, \Delta Y_s = 0} V_{*s\vartheta\tau}(\Delta X_s) \exp \left(- \int_0^t K_{*s\vartheta} (V_{*s\vartheta\tau} - 1) ds \right) \\ & \prod_{s \leq t: \Delta Y_s \neq 0} V_{-*s\vartheta\tau}(\Delta X_s, \Delta Y_s). \end{aligned}$$

For a heuristic derivation in terms of product integrals, see Andersen et al. (1993, p. 107). We have already noted in Remark 5 that the derivative of the first factor, the partial likelihood ratio, equals the partial score function

(3.10). The derivative of the second factor is obtained similarly. Finally, $\partial_{\tau=\vartheta} V_{-*s\vartheta\tau} = \ell'_{-*s\vartheta}$ by definition.

The representation (3.13) of the likelihood can be used in the partially specified model of Result 2 to prove that the optimal estimating function obtained there is efficient as long as no additional restrictions involving ϑ are imposed on the model. The arguments are similar to those outlined in Section 1 for the i.i.d. case. In Greenwood and Wefelmeyer (1996) a representation analogous to (3.13) is given for general semimartingales and can be used to generalize the results obtained here to partially specified semimartingale regression models.

References

- Andersen, P. K., Borgan, Ø., Gill, R. D. and Keiding, N. (1993). *Statistical Models Based on Counting Processes*. Springer-Verlag, New York.
- Arjas, E. and Haara, P. (1984). A marked point process approach to censored failure data with complicated covariates. *Scand. J. Statist.* 11, 193–209.
- Bickel, P.J., Klaassen, C.A.J., Ritov, Y. and Wellner, J.A. (1993). *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins University Press, Baltimore.
- Cox, D. R. (1975). Partial likelihood. *Biometrika* 62, 269–276.
- Greenwood, P. E. and Wefelmeyer, W. (1996). Cox's factoring of regression model likelihoods for continuous time processes. To appear in *Bernoulli*.
- Heyde, C.C. (1987). On combining quasi-likelihood estimating functions. *Stochastic Process. Appl.* 25 281–287.
- Jacod, J. (1987). Partial likelihood process and asymptotic normality. *Stochastic Process. Appl.* 26, 47–71.
- Jacod, J. (1990). Sur le processus de vraisemblance partielle. *Ann. Inst. Henri Poincaré* 26, 299–329.
- Jacod, J. and Shiryaev, A. N. (1987). *Limit Theorems for Stochastic Processes*. Grundlehren der mathematischen Wissenschaften 288, Springer-Verlag, Berlin.

