

Optimality and small Δ -optimality of martingale estimating functions

Martin Jacobsen*

Abstract

Martingale estimating functions determined from a given collection (the base) of conditional expectations are considered for estimating the parameters of a discretely observed diffusion. It is discussed how to make the martingale estimating functions small Δ -optimal, i.e. nearly efficient when the observations are close together, in particular it is shown that this is possible provided the base is large enough. It is also shown that the optimal martingale estimating function with a given base, is automatically small Δ -optimal, provided only that the base is sufficiently large. In both cases the critical dimension of the base is the same and determined by the dimension of the diffusion, and on whether the squared diffusion matrix is parameter dependent or not; the critical number does not depend however on the dimension of the parameter.

KEYWORDS AND PHRASES: unbiased estimating functions, minimizing asymptotic covariances, conditional moments, a generalized Cox-Ingersoll-Ross process, Gaussian diffusions.

1 Introduction

Suppose a d -dimensional, ergodic, timehomogeneous diffusion X is observed at finitely many points $i\Delta$ in time, $i = 0, \dots, n$. In order to estimate the unknown parameter θ that determines the distribution of X , rather than doing maximum-likelihood, which may well prove unfeasible, one often resorts to

*MaPhySto – Centre for Mathematical Physics and Stochastics, funded by a grant from the Danish National Research Foundation.

the use of unbiased estimating functions, some of the most successful of which are based on conditional expectations, resulting in estimating equations of the form

$$\sum_{i=1}^n \sum_{q=1}^r h^q(X_{(i-1)\Delta}) (f^q(X_{i\Delta}) - E_\theta(f^q(X_{i\Delta}) | X_{(i-1)\Delta})) = 0 \quad (1)$$

where, for the moment, we consider θ to be one-dimensional (see (6) and (2) below for the general setup). (1) is the prime example of an estimating equation obtained from an unbiased *martingale* estimating function.

The study of estimating equations of the form (1) was initiated by Bibby and Sørensen [1] who focused on the case $r = 1$ and $f(x) = x$ and for that case also showed that under mild conditions (1) has a consistent and \sqrt{n} -asymptotically Gaussian solution in θ (as $n \rightarrow \infty$ with $\Delta > 0$ fixed). These asymptotic results readily generalize to other types of unbiased estimating equations, see e.g. Sørensen [6].

In the present paper we consider estimating equations of the form (1) for d -dimensional diffusions X with a p -dimensional parameter θ . The main issue is the discussion of the choice of *base* $(f^q)_{1 \leq q \leq r}$, in particular the choice of r , and the choice of *weights* h^q . For this we think of $\Delta > 0$ as arbitrary and then consider families of estimating equations of the form (1) where h^q , but not f^q , is allowed to depend on Δ , and such that for any Δ , (1) has a consistent and asymptotically Gaussian solution as described above.

Given the base (f^q) there is an *optimal* choice of weights, see Proposition 1 below. The resulting estimator will typically not be efficient, while an estimator that is *small Δ -optimal* will be nearly efficient for small values of Δ – and of course still consistent and asymptotically Gaussian for all Δ , although not optimal. As we shall see, while it may be difficult to find the optimal estimator (to obtain and use the weights one needs the explicit form of the inverse of a $r \times r$ -matrix with all elements conditional variances), it is quite easy to determine weights that lead to small Δ -optimal estimators.

The concept of small Δ -optimality was introduced by Jacobsen [2], and the main purpose of the present paper is to discuss conditions for small Δ -optimality for the type of martingale estimating functions underlying (1).

The first main result, Theorem 2, shows that given the base, provided only that the dimension r is large enough, there always exist weights such that small Δ -optimality is achieved. Furthermore it is easy to find the

weights – it is just a matter of solving at any point x in the range for the diffusion, a set of linear equations, and in the statement of the theorem, a concrete solution is exhibited.

The second main result, Theorem 3, shows that for any base, again provided r is large enough, for the optimal choice of weights small Δ –optimality is automatic.

In both theorems the same critical value r_0 for the dimension of the base appears: for $r \geq r_0$ small Δ –optimality can be achieved for any base, while for $r < r_0$ this may only be possible, if at all, for a special choice of base. (The case $r < r_0$ is only mentioned, but not really explored below). The value of r_0 depends on the structure of the model but not on the dimension of the parameter with $r_0 = d$, the dimension of the diffusion, if the diffusion coefficient does not depend on the parameter, and $r_0 = d(d+3)/2$ otherwise. Thus one natural choice of base is $f^i(x_1, \dots, x_d) = x_i$ if $r_0 = d$, and these f^i supplemented by $f^{ij}(x_1, \dots, x_d) = x_i x_j$ for $1 \leq i \leq j \leq d$ if $r_0 = d(d+3)/2$.

The paper is concluded (Section 3) with two examples, one describing a generalized Cox-Ingersoll-Ross model, the second the finite-dimensional Ornstein-Uhlenbeck processes. For the latter it turns out, that using the base of first and second order moments (f^i, f^{ij}) , the concrete small Δ –optimal estimating function exhibited in Theorem 2 for any Δ yields the maximum-likelihood estimator.

Although both here and in Jacobsen [2], small Δ –optimality is discussed exclusively for diffusions, we wish to point out that the concept makes perfect sense for any model involving discrete observations from an ergodic, timehomogeneous Markov process in continuous time.

2 Optimality and small Δ –optimality

Let $X = (X_t)_{t \geq 0}$ be a d –dimensional ergodic diffusion, solving the stochastic differential equation

$$dX_t = b_\theta(X_t) dt + \sigma_\theta(X_t) dB_t, \quad X_0 = U$$

where $b_\theta(x) \in \mathbb{R}^{d \times 1}$, $\sigma_\theta(x) \in \mathbb{R}^{d \times d}$, B is a standard d –dimensional Brownian motion and U is a d –dimensional random variable, independent of B . Both the drift b_θ and the diffusion coefficient σ_θ is allowed to depend on the p –dimensional parameter $\theta \in \Theta$. The invariant distribution for X is de-

noted μ_θ , i.e. if U has distribution μ_θ , then X is a strictly stationary Markov process.

We shall also assume that X takes its values within some open subset D of \mathbb{R}^d , of course not allowed to depend on θ . Also we assume that $C_\theta(x) := \sigma_\theta(x)\sigma_\theta^T(x) \succ 0$ for all θ and all $x \in D$, i.e. the symmetric positive semidefinite matrix $C_\theta(x)$ is assumed to be strictly positive definite always. (T denotes matrix transposition).

We shall write μ_θ also for the density of μ_θ and assume that for all θ , $\mu_\theta > 0$ everywhere on D . The transition density is denoted $p_{t,\theta}(x, y)$,

$$P_\theta(X_{s+t} \in dy | X_s = x) = p_{t,\theta}(x, y) dy.$$

The underlying probability P_θ depends not only on θ but also on the distribution of X_0 . It is denoted P_θ^μ if X_0 has distribution μ_θ and P_θ^x if $X_0 \equiv x$. The corresponding expectations are written E_θ^μ and E_θ^x .

We shall denote by $Q_{t,\theta}$ the joint distribution of (X_s, X_{s+t}) under P_θ^μ (for any s), and by $q_{t,\theta}$ the density of $Q_{t,\theta}$,

$$q_{t,\theta}(x, y) = \mu_\theta(x)p_{t,\theta}(x, y).$$

Finally, the transition operators for X are denoted $\pi_{t,\theta}$,

$$\pi_{t,\theta}f(x) = E_\theta^x f(X_t)$$

provided the integral makes sense (e.g. for f bounded or $f \in L^1(\mu_\theta)$), and the differential operator determining the infinitesimal generator is denoted A_θ ,

$$A_\theta f(x) = \sum_{i=1}^d b_\theta^i(x) \partial_{x_i} f(x) + \frac{1}{2} \sum_{i,j=1}^d C_\theta^{ij}(x) \partial_{x_i x_j}^2 f(x)$$

for sufficiently smooth functions f .

Suppose now that X is observed at finitely many equidistant timepoints, $X_0, X_\Delta, \dots, X_{n\Delta}$. We shall discuss optimality properties of estimators based on martingale estimating functions, i.e. the estimator $\hat{\theta}_n$ of θ is found by solving the estimating equation

$$G_{n,\Delta}(\theta) = \sum_{i=1}^n g_{\Delta,\theta}(X_{(i-1)\Delta}, X_{i\Delta}) = 0 \tag{2}$$

where $(t, \theta, x, y) \rightarrow g_{t,\theta}(x, y)$ is a p -variate function such that each coordinate $g_{t,\theta}^k$ satisfies the martingale condition

$$g_{t,\theta}^{k*}(x) = 0 \text{ for all } x, \quad g_{t,\theta}^{k*}(x) = E_{\theta}^x g_{t,\theta}^k(X_0, X_t) \quad (3)$$

ensuring that $(G_{n,\Delta}(\theta))_{n \geq 1}$ is a p -dimensional P_{θ} -martingale (whatever the initial distribution of X).

Following the terminology in Jacobsen [2], we shall refer to $\mathcal{G} = (g_{t,\theta})_{t>0, \theta \in \Theta}$ as a *well behaved flow of martingale estimating functions*, $\mathcal{G} \subset \mathcal{M}$ (the space of flows of martingale estimating functions) if each $g_{t,\theta}^k \in L^2(Q_{t,\theta})$ with

$$E_{\theta}^{\mu} g_{t,\theta'}(X_0, X_t) = 0 \text{ if and only if } \theta = \theta'$$

(where the equality holds for $\theta = \theta'$ because of (3)) and if furthermore $E_{\theta}^{\mu}(g_{t,\theta} g_{t,\theta}^T)(X_0, X_t) \succ 0$ for all t, θ , $\partial_{\theta_t} g_{t,\theta}^k \in L^1(Q_{t,\theta})$ for all t, θ , all indices k, ℓ , and finally and most important if for every θ_0 and every $t = \Delta > 0$ there is with $P_{\theta_0}^{\mu}$ -probability tending to 1 a consistent solution $\hat{\theta}_n$ to (2) such that $\sqrt{n}(\hat{\theta}_n - \theta_0)$ converges in distribution for $n \rightarrow \infty$ to the p -dimensional Gaussian distribution with mean vector 0 and covariance matrix $\text{var}_{\Delta, \theta_0}(g, \hat{\theta})$ given by

$$\text{var}_{\Delta, \theta_0}(g, \hat{\theta}) = \Lambda_{\Delta, \theta_0}^{-1}(g) E_{\theta_0}^{\mu}(g_{\Delta, \theta_0} g_{\Delta, \theta_0}^T)(X_0, X_{\Delta}) (\Lambda_{\Delta, \theta_0}^{-1}(g))^T. \quad (4)$$

Here

$$\Lambda_{t,\theta} := E_{\theta}^{\mu} \dot{g}_{t,\theta}(X_0, X_t), \quad (5)$$

the dot $\dot{\cdot}$ signifying differentiation with respect to θ so that $\dot{g}_{t,\theta}(x, y) \in \mathbb{R}^{p \times p}$ is given by

$$(\dot{g}_{t,\theta}(x, y))_{k\ell} = \partial_{\theta_t} g_{t,\theta}^k(x, y).$$

The reader is reminded that asymptotic normality of $\hat{\theta}_n$, as specified above, holds under quite weak assumptions (see Sørensen [6]) and that certainly (4) is the natural expression for the asymptotic covariance. The most critical among the assumptions needed is that $\Lambda_{t,\theta} \in \mathbb{R}^{p \times p}$ must be non-singular for all t, θ .

Notation. Throughout the paper, derivatives are understood as matrices in analogy with (5): if ϕ is a ρ -variate function of a v -dimensional variable $z = (z_1, \dots, z_v) \in \mathbb{R}^v$, $\partial_z \phi$ denotes the $\rho \times v$ -matrix of partial derivatives

with κ 'th row $(\partial_{z_1} \phi^\kappa, \dots, \partial_{z_v} \phi^\kappa)$. The dot notation is used exclusively for differentiation with respect to θ , $\dot{\phi} = \partial_\theta \phi$.

In the remainder of the paper we shall focus on martingale estimating functions derived from conditional expectations of given functionals, i.e. we assume that

$$g_{t,\theta}^k(x, y) = \sum_{q=1}^r h_{t,\theta}^{qk}(x) (f_\theta^q(y) - (\pi_{t,\theta} f_\theta^q)(x)), \quad (6)$$

or, in matrix notation,

$$g_{t,\theta}(x, y) = h_{t,\theta}^T(x) (f_\theta(y) - (\pi_{t,\theta} f_\theta)(x))$$

with $h_{t,\theta}(x) \in \mathbb{R}^{r \times p}$, $f_\theta(x) \in \mathbb{R}^{r \times 1}$. (The integrability assumptions imposed on general g above, makes it natural to assume here that f_θ^q and $h_{t,\theta}^{qk} \in L^4(\mu_\theta)$, while θ -derivatives of f_θ^q and $h_{t,\theta}^{qk}$ must belong to $L^2(\mu_\theta)$. We shall not be too concerned about these conditions in the sequel – it is tacitly assumed everywhere that the flow \mathcal{G} given by (6) is well behaved).

Estimating functions of the form (6) were first used by Bibby and Sørensen [1], see also Jacobsen [2], Section 3 for an overview.

We shall refer to the functions $f_\theta^1, \dots, f_\theta^r$ as the *base* for the flow of estimating functions given by (6). The problem to be studied in this paper is that of finding good choices for the dimension of the base and for the *weights* $h_{t,\theta}$ given the base.

Assumption A. *The functions $f_\theta^q(x)$ are supposed to be differentiable in θ and twice differentiable in x . Also, the base $f_\theta^1, \dots, f_\theta^r$ is supposed to have full affine rank r on the domain D for all θ , i.e. for an arbitrary θ the identity*

$$\sum_{q=1}^r a_\theta^q f_\theta^q(x) + \alpha_\theta = 0 \quad (x \in D)$$

for some constants $a_\theta^q, \alpha_\theta$ implies $a_\theta^1 = \dots = a_\theta^r = \alpha_\theta = 0$.

The functions $h_{t,\theta}^{qk}$ are supposed to satisfy that for any t, θ , the p r -variate functions $x \rightarrow (h_{t,\theta}^{1k}(x), \dots, h_{t,\theta}^{rk}(x))$ forming the columns of $h_{t,\theta}$ are linearly independent on D . ■

Note that if $f_\theta^1, \dots, f_\theta^r$ does not have full affine rank, there is a representation (6) of the $g_{t,\theta}^k$ with r replaced by $r - 1$. The condition that $f_\theta^1, \dots, f_\theta^r$ has full rank is equivalent to assuming that the r d -variate functions $\partial_x f_\theta^q$

for $1 \leq q \leq r$ be linearly independent. In the main results, Theorems 2 and 3 below, Assumption A is supplemented by conditions on the pointwise behaviour of $\partial_x f_\theta$ and $\partial_{xx}^2 f_\theta$.

If for some t it holds for all θ that the columns of $h_{t,\theta}$ are not linearly independent, i.e there is $\beta_{t,\theta} \in \mathbb{R}^{p \times 1} \setminus 0$ such that $h_{t,\theta}(x)\beta_{t,\theta} = 0$ for all x , then $\beta_{t,\theta}^T g_{t,\theta}(x, y) = 0$ for all x, y so that one of the p estimating equations in (2) can be obtained from the others and it is impossible to estimate all p parameters θ_ℓ – formally both matrices $\Lambda_{t,\theta}(g)$ and $E_\theta^\mu(g_{t,\theta} g_{t,\theta}^T)(X_0, X_t)$ become singular and (4) does not make sense.

Note that we allow the base $f_\theta^1, \dots, f_\theta^r$ to depend on θ , but *not* on t .

For a given base, it is easy to determine the optimal choices for the $h_{t,\theta}^{qk}$, i.e. the choices minimizing $\text{var}_{t,\theta}(g, \hat{\theta})$. We use the notation $A \succeq B$ between symmetric, positive semidefinite matrices to signify that $A - B$ is also positive semidefinite.

Proposition 1 *Assume that for all x, t and θ the symmetric $r \times r$ –matrix*

$$\Pi_{t,\theta} f_\theta(x) := \pi_{t,\theta}(f_\theta f_\theta^T)(x) - (\pi_{t,\theta} f_\theta)(x) (\pi_{t,\theta} f_\theta^T)(x)$$

is non-singular, and define

$$h_{t,\theta}^{\text{opt}} = (\Pi_{t,\theta} f_\theta)^{-1} \left(\partial_\theta (\pi_{t,\theta} f_\theta) - \pi_{t,\theta} \left(\dot{f}_\theta \right) \right). \quad (7)$$

Then, provided differentiation and integration can be interchanged in

$$\partial_\theta \int p_{t,\theta}(x, y) f_\theta(y) dy = \int \partial_\theta (p_{t,\theta}(x, y) f_\theta(y)) dy$$

and the flow \mathcal{G}^{opt} given by

$$g_{t,\theta}^{\text{opt}}(x, y) = (h_{t,\theta}^{\text{opt}})^T(x) (f_\theta(y) - \pi_{t,\theta} f_\theta(x)) \quad (8)$$

is well behaved, it holds that

$$\text{var}_{t,\theta}(g^{\text{opt}}, \hat{\theta}) \preceq \text{var}_{t,\theta}(g, \hat{\theta})$$

for any well behaved flow $\mathcal{G} = (g_{t,\theta})$ of the form (6) with base $f_\theta^1, \dots, f_\theta^r$.

Proof. Since f is allowed to depend on θ , this extends (2.10) in Bibby and Sørensen [1] and Example 5.5 in Jacobsen [2], so we indicate the proof. By the projection theorem (Kessler [3], Proposition 1, Jacobsen [2], Proposition 5.3), $g_{t,\theta}^{k,\text{opt}}$ is found by projecting the k 'th coordinate of the score function, $\partial_{\theta_k} p_{t,\theta}(x, y)/p_{t,\theta}(x, y)$ onto the subspace of $L^2(Q_{t,\theta})$ spanned by functions of the form (6) with the f_θ^q fixed and arbitrary $h_{t,\theta}^{qk} \in L^2(\mu_\theta)$. Thus $h_{t,\theta}^{qk,\text{opt}}$ satisfies for all $1 \leq q_0 \leq r$, $1 \leq k \leq p$ and all $h \in L^2(\mu_\theta)$ that

$$0 = E_\theta^\mu \left[\frac{\partial_{\theta_k} p_{t,\theta}}{p_{t,\theta}}(X_0, X_t) - \sum_{q=1}^r h_{t,\theta}^{qk,\text{opt}}(X_0) (f_\theta^q(X_t) - \pi_{t,\theta} f_\theta^q(X_0)) \right] \\ \times h(X_0) (f_\theta^{q_0}(X_t) - \pi_{t,\theta} f_\theta^{q_0}(X_0)) \quad (9)$$

Using that $\dot{p}_{t,\theta}/p_{t,\theta}$ is a martingale estimating function and that

$$E_\theta^x \frac{\partial_{\theta_k} p_{t,\theta}}{p_{t,\theta}}(X_0, X_t) f_\theta^{q_0}(X_t) = \int \partial_{\theta_k} p_{t,\theta}(x, y) f_\theta^{q_0}(y) dy \\ = \partial_{\theta_k} \int p_{t,\theta}(x, y) f_\theta^{q_0}(y) dy \\ - \int p_{t,\theta}(x, y) \partial_{\theta_k} f_\theta^{q_0}(y) dy,$$

(9) may be written

$$E_\theta^\mu h(X_0) [\partial_{\theta_k} (\pi_{t,\theta} f_\theta^{q_0})(X_0) - \pi_{t,\theta} (\partial_{\theta_k} f_\theta^{q_0})(X_0) \\ - \sum_{q=1}^r h_{t,\theta}^{qk,\text{opt}}(X_0) (\pi_{t,\theta} (f_\theta^q f_\theta^{q_0})(X_0) - \pi_{t,\theta} f_\theta^q(X_0) \pi_{t,\theta} f_\theta^{q_0}(X_0))]]$$

for all $h \in L^2(\mu_\theta)$, i.e. the expression in square brackets must vanish P_θ^μ -a.s. and the result follows. \blacksquare

Proposition 1 is a result on *local optimality*, i.e. it exhibits the best member from a given, restricted class of estimating functions, best from the point of view of minimizing the asymptotic covariance of the resulting estimator. But only in exceptional cases will this choice be *globally optimal*, i.e. the (locally) optimal estimator will be efficient against the maximum-likelihood estimator.

By contrast, the concept of small Δ -optimality introduced by Jacobsen [2], Section 7, gives conditions for global optimality, not for any given $\Delta > 0$,

but only for $\Delta \rightarrow 0$. We shall briefly recapitulate the sufficient conditions for small Δ -optimality of martingale estimating functions, Theorem 7.5 in Jacobsen [2].

With $\mathcal{G} \subset \mathcal{M}$ a well behaved flow of estimating functions, it is first of all essential to assume that there is a smooth extension of $g_{t,\theta}(x, y)$ (which is defined only for $t > 0$) to allow $t = 0$, i.e. after a possible renormalization of $g_{t,\theta}$ by a factor (non-zero scalar or non-singular $p \times p$ -matrix) depending on t, θ but not on x, y (so the solution of (2) is not affected) the limit

$$g_{0,\theta}(x, y) = \lim_{t \rightarrow 0} g_{t,\theta}(x, y)$$

must exist with $(x, y) \rightarrow g_{0,\theta}(x, y)$ not identically 0, of full rank p in a suitable sense and sufficiently smooth as required by the conditions below.

With this smooth extension of $g_{t,\theta}$ available, it is shown in Jacobsen [2] that subject to important integrability conditions (see the paragraph below (13)), the asymptotic covariance for $\hat{\theta}_n$ has an expansion, as $\Delta \rightarrow 0$,

$$\text{var}_{\Delta,\theta} \left(g, \hat{\theta} \right) = \frac{1}{\Delta} v_{-1,\theta} \left(g, \hat{\theta} \right) + v_{0,\theta} \left(g, \hat{\theta} \right) + o(1) \quad (10)$$

and three cases (i), (ii) and (iii) for the structure of the diffusion model are then considered for the discussion of small Δ -optimality (to achieve the structure in (iii) it may be necessary first to reparametrize the model):

- (i) $C_\theta = C$ does not depend on θ . Then the main term in (10) is always present and small Δ -optimality is achieved by minimizing globally (over all g) $v_{-1,\theta} \left(g, \hat{\theta} \right)$. A sufficient condition for a given flow $(g_{t,\theta})$ to be small Δ -optimal is that

$$\partial_y g_{0,\theta}(x, x) = K_\theta \dot{b}_\theta^T(x) C^{-1}(x) \quad (11)$$

for some non-singular $K_\theta \in \mathbb{R}^{p \times p}$. ($\partial_y g_{0,\theta}(x, x)$ evaluates $\partial_y g_{0,\theta}(x, y)$ along the diagonal $y = x$).

- (ii) C_θ depends on all parameters $\theta_1, \dots, \theta_p$. Then the main term in (10) vanishes provided $\partial_y g_{0,\theta}(x, x) \equiv 0$ and small Δ -optimality is achieved by minimizing $v_{0,\theta} \left(g, \hat{\theta} \right)$. A sufficient condition for $(g_{t,\theta})$ to be small Δ -optimal is that

$$\partial_y g_{0,\theta}(x, x) = 0, \quad \partial_{yy}^2 g_{0,\theta}(x, x) = K_\theta \dot{C}_\theta^T(x) \left(C_\theta^{\otimes 2}(x) \right)^{-1} \quad (12)$$

for some non-singular $K_\theta \in \mathbb{R}^{p \times p}$. (Here $\dot{C}_\theta(x) \in \mathbb{R}^{d^2 \times p}$ with $(\dot{C}_\theta(x))_{ij,k} = \partial_{\theta_k} C_\theta^{ij}(x)$).

- (iii) C_θ depends on the parameters $\theta_1, \dots, \theta_{p'}$ but not on $\theta_{p'+1}, \dots, \theta_p$ for some p' with $1 \leq p' < p$. Then parts of the main term in (10) can be made to disappear so that

$$v_{-1,\theta}(g, \hat{\theta}) = \begin{pmatrix} 0_{p' \times p'} & 0_{p' \times (p-p')} \\ 0_{(p-p') \times p'} & v_{22,-1,\theta}(g, \hat{\theta}) \end{pmatrix}.$$

Furthermore the matrix $v_{22,-1,\theta}(g, \hat{\theta}) \in \mathbb{R}^{(p-p') \times (p-p')}$ can be minimized and small Δ -optimality is achieved by in addition minimizing the upper left block $v_{11,0,\theta}(g, \hat{\theta})$ of $v_{0,\theta}(g, \hat{\theta})$. A sufficient condition for small Δ -optimality is that

$$\begin{aligned} \partial_y g_{0,\theta}(x, x) &= c_\theta \begin{pmatrix} 0_{p' \times d} \\ \dot{b}_{2,\theta}^T(x) C_\theta^{-1}(x) \end{pmatrix} \\ \partial_{yy}^2 g_{1,0,\theta}(x, x) &= K'_\theta \dot{C}_{1,\theta}^T(x) (C_\theta^{\otimes 2}(x))^{-1} \end{aligned} \quad (13)$$

for some constant $c_\theta \neq 0$ and some non-singular $K'_\theta \in \mathbb{R}^{p' \times p'}$. (Notation: $\dot{b}_{2,\theta} \in \mathbb{R}^{d \times (p-p')}$ comprises the last $p-p'$ columns of \dot{b}_θ , $g_{1,0,\theta}$ the first p' coordinates of $g_{0,\theta}$ and $\dot{C}_{1,\theta} \in \mathbb{R}^{d^2 \times p'}$ the first p' columns of \dot{C}_θ).

As mentioned above, to check for small Δ -optimality more is required than just checking (11), (12) or (13), viz. it must be verified that various matrices involving expectations of quantities related to \dot{b} , \dot{C} , $\partial_y g_{0,\theta}$ and $\partial_{yy}^2 g_{0,\theta}$ must be non-singular, see Theorem 7.5 in Jacobsen [2].

Remark 1 *The conditions on $\partial_y g_{0,\theta}$ in (12) and (13) are important: if not satisfied, the main term in (10) will not disappear and some of the components of the estimates resulting from g , will be totally inefficient – have efficiency close to 0 – if Δ is small.*

We shall now show that subject to these integrability conditions, small Δ -optimality of martingale estimating functions is easy to achieve. The three cases refer to (i), (ii) and (iii) above.

Notation. Let $J := \{(i', j') : 1 \leq i' \leq j' \leq d\}$. Thus J has $|J| = d + \binom{d}{2}$ elements and can be used as an index set for characterizing the elements of

a symmetric $d \times d$ -matrix. We write $R \in \mathbb{R}^{d^2 \times J}$ for the reduction matrix with elements

$$R_{ij,i'j'} = \delta_{ii'}\delta_{jj'} \quad (1 \leq i, j \leq d, (i', j') \in J).$$

Thus, if $M \in \mathbb{R}^{A \times d^2}$, $MR \in \mathbb{R}^{A \times J}$ with $(MR)_{a,i'j'} = M_{a,i'j'}$ as is used frequently below.

As counterpart to R , the expansion matrix $\tilde{R} \in \mathbb{R}^{J \times d^2}$ is defined by

$$\tilde{R}_{i'j',ij} = \begin{cases} \delta_{i'i}\delta_{j'j} & \text{if } i \leq j \\ \delta_{i'j}\delta_{j'i} & \text{if } i > j. \end{cases}$$

Then

$$\tilde{R}R = I_{J \times J} \tag{14}$$

and for any matrix $N \in \mathbb{R}^{S \times d^2}$, symmetric in the sense that $N_{s,ij} = N_{s,ji}$ for all, s, i, j , it holds that

$$N(R\tilde{R}) = N. \tag{15}$$

Define

$$\dim(d) := d + |J| = d(d + 3)/2,$$

a number that plays a critical role below.

Theorem 2 *Let $(f_\theta^1, \dots, f_\theta^r)$ be a base for a martingale estimating function, of full affine rank r .*

- (i) *Suppose that $r \geq d$, that for μ_θ -a.a. x the matrix $\partial_x f_\theta(x) \in \mathbb{R}^{r \times d}$ is of full rank d , and that the p d -variate functions forming the columns of b_θ are linearly independent. Then there exists $h_{t,\theta}(x) = h_\theta(x) \in \mathbb{R}^{r \times p}$ not depending on t such that $g_{t,\theta}(x, y) := h_\theta^T(x) (f_\theta(y) - \pi_{t,\theta} f(x))$ satisfies the small Δ -optimality condition (11). In particular, for $r = d$ one may choose*

$$h_\theta^T(x) = b_\theta^T(x) C^{-1}(x) (\partial_x f_\theta(x))^{-1} \tag{16}$$

and this h_θ has linearly independent columns as required in Assumption A.

- (ii) *Suppose that $r \geq \dim(d)$, that for μ_θ -a.a. x , the matrix*

$$\begin{pmatrix} \partial_x f_\theta(x) & \partial_{xx}^2 f_\theta(x) \end{pmatrix} \in \mathbb{R}^{r \times (d+d^2)}$$

is of full rank $\dim(d)$ and that the p d^2 -variate functions forming the columns of \dot{C}_θ are linearly independent. Then there exists $h_{t,\theta} = h_\theta \in \mathbb{R}^{r \times p}$ not depending on t such that $g_{t,\theta}(x, y) := h_\theta^T(x) (f_\theta(y) - \pi_{t,\theta} f(x))$ satisfies the small Δ -optimality condition (12). In particular, for $r = \dim(d)$ one may choose

$$h_\theta^T(x) = \left(\begin{array}{cc} 0_{p \times d} & \dot{C}_\theta^T(x) (C_\theta^{\otimes 2}(x))^{-1} R \end{array} \right) \left(\begin{array}{cc} \partial_x f_\theta(x) & \partial_{xx}^2 f_\theta(x) R \end{array} \right)^{-1}, \quad (17)$$

and this h_θ has linearly independent columns as required in Assumption A.

(iii) Suppose that $r \geq \dim(d)$, that for μ_θ -a.a. x , the matrix

$$\left(\begin{array}{cc} \partial_x f_\theta(x) & \partial_{xx}^2 f_\theta(x) \end{array} \right) \in \mathbb{R}^{r \times (d+d^2)}$$

is of full rank $\dim(d)$, that the $p - p'$ d -variate functions forming the columns of $\dot{b}_{2,\theta}$ are linearly independent, and that the p' d^2 -variate functions forming the columns of $\dot{C}_{1,\theta}$ are linearly independent. Then there exists $h_{t,\theta} = h_\theta \in \mathbb{R}^{r \times p}$ not depending on t such that $g_{t,\theta}(x, y) := h_\theta^T(x) (f_\theta(y) - \pi_{t,\theta} f(x))$ satisfies the small Δ -optimality condition (13). In particular, for $r = \dim(d)$ one may choose

$$h_\theta^T(x) = \left(\begin{array}{cc} 0_{p' \times d} & \dot{C}_{1,\theta}^T(x) (C_\theta^{\otimes 2}(x))^{-1} R \\ \dot{b}_{2,\theta}^T(x) C_\theta^{-1}(x) & * \end{array} \right) \left(\begin{array}{cc} \partial_x f_\theta(x) & \partial_{xx}^2 f_\theta(x) R \end{array} \right)^{-1}. \quad (18)$$

with $*$ a $(p - p') \times |J|$ matrix, depending arbitrarily on θ and x . If $*$ is chosen $= 0$, then this h_θ has linearly independent columns as required in Assumption A.

Proof. Since h_θ does not depend on t ,

$$g_{0,\theta}(x, y) = h_\theta^T(x) (f_\theta(y) - f_\theta(x))$$

whence

$$\partial_y g_{0,\theta}(x, x) = h_\theta^T(x) \partial_x f_\theta(x), \quad \partial_{yy}^2 g_{0,\theta}(x, x) = h_\theta^T(x) \partial_{xx}^2 f_\theta(x).$$

Thus, for each x , (11), (12) or (13) gives a system of linear equations for determining the elements of $h_\theta(x)$. The conditions of the theorem ensures

that these equations have at least one solution, and exactly one in case (i) if $r = d$ and in case (ii) if $r = \dim(d)$. (For case (ii), note that since $\partial_{x_i x_j}^2 = \partial_{x_j x_i}^2$, the rank of $\partial_{xx}^2 f_\theta$ is at most $|J|$. With h_θ given by (17), one now finds

$$\partial_{yy}^2 g_{0,\theta}(x, x)R = \dot{C}_\theta^T(x) (C_\theta^{\otimes 2}(x))^{-1} R$$

and using (15) this implies the second identity in (12).

The assertions about h_θ having linearly independent columns follow readily from the assumptions made on the columns of \dot{b}_θ (case (i)), \dot{C}_θ (case (ii)) and $\dot{b}_{2,\theta}$ and $\dot{C}_{1,\theta}$ (case (iii)). \blacksquare

Theorem 2 only gives a solution for h_θ such that the relevant of (11), (12) or (13) is satisfied. To check small Δ -optimality one further has to check the required integrability conditions (e.g. that all $h_\theta^{qk} \in L^2(\mu_\theta)$).

In Theorem 2 we have exhibited a concrete choice of small Δ -optimal estimating functions from a given base (f_θ^q) . But it is then easy to define a host of others that are also small Δ -optimal, but may behave better for a given Δ , viz. flows $(g_{t,\theta})$ of the form

$$g_{t,\theta}^k(x, y) = \sum_{q=1}^r a_\theta^{qk}(t) h_\theta^{qk}(x) (f_\theta^q(y) - \pi_{t,\theta} f_\theta^q(x)) \quad (19)$$

with h_θ^T given by the relevant of (16), (17) or (18) and each $a_\theta^{qk}(t)$ a non-random function of t , continuous with $a_\theta^{qk}(0) = 1$: for this flow, $g_{0,\theta}$ is the same as for the original flow, so small Δ -optimality still holds. However, there is no obvious optimal choice for the $a_\theta^{qk}(t)$: each $g_{t,\theta}^k$ given by (19) varies in a linear subspace when the $a_\theta^{qk}(t)$ are arbitrary, but the subspace depends on k so the projection technique from the proof of Proposition 1 does not apply. (To use the proposition and find optimal constants, one must consider a much larger space of estimating functions such as

$$g_{t,\theta}^k(x, y) = \sum_{k'=1}^p \sum_{q=1}^r a_\theta^{k,qk'}(t) h_\theta^{qk'}(x) (f_\theta^q(y) - \pi_{t,\theta} f_\theta^q(x))$$

with arbitrary constants $a_\theta^{k,qk'}(t)$. The optimal constants can be found explicitly, but the expression involves the for practical purposes unwieldy inverse of the $pr \times pr$ -matrix with $(kq, k'q')$ 'th element

$$E_\theta^\mu h_\theta^{qk}(X_0) (f_\theta^q(X_t) - \pi_{t,\theta} f_\theta^q(X_0)) h_\theta^{q'k'}(X_0) \left(f_\theta^{q'}(X_t) - \pi_{t,\theta} f_\theta^{q'}(X_0) \right).$$

Remark 2 In Theorem 2, case (iii), the expression (18) for $h_\theta^T(x)$ depends on the choice of $*$. A different small Δ -optimal flow not involving this arbitrary matrix may be obtained as follows: take $r = \dim(d)$ and apart from the given base (f_θ^q) of dimension r , choose a second base (\tilde{f}_θ^q) of dimension d . Assuming in addition to the assumptions from the theorem that $\partial_x \tilde{f}_\theta(x)$ is non-singular, the flow $(\tilde{g}_{t,\theta})$ given by the p' , respectively $p - p'$ components

$$\begin{aligned}\tilde{g}_{1,t,\theta}(x, y) &= h_{1,\theta}^T(x) (f_\theta(y) - \pi_{t,\theta} f_\theta(x)) \\ \tilde{g}_{2,t,\theta}(x, y) &= \tilde{h}_\theta(x) (\tilde{f}_\theta(y) - \pi_{t,\theta} \tilde{f}_\theta(x))\end{aligned}$$

with

$$\begin{aligned}h_{1,\theta}^T(x) &= \left(0_{p' \times d} \quad \dot{C}_{1,\theta}^T(x) (C_{1,\theta}^{\otimes 2}(x))^{-1} R \right) \left(\partial_x f_\theta(x) \quad \partial_{xx}^2 f_\theta(x) R \right)^{-1} \\ \tilde{h}_\theta^T(x) &= \tilde{b}_{2,\theta}^T(x) C_\theta^{-1}(x) \left(\partial_x \tilde{f}_\theta(x) \right)^{-1}\end{aligned}\tag{21}$$

is small Δ -optimal according to (13).

For general flows $(\tilde{g}_{t,\theta})$ of the form (21), with the components of $\tilde{g}_{1,t,\theta}$ and $\tilde{g}_{2,t,\theta}$ varying in two different subspaces, optimal time dependent choices $h_{1,t,\theta}^T$ of $h_{1,\theta}^T$ and $\tilde{h}_{t,\theta}^T$ of \tilde{h}_θ^T cannot be found using e.g. Proposition 1. However partly optimal candidates may be determined using the proposition twice: first with base (f_θ^q) in the model with $\theta_{p'+1}, \dots, \theta_p$ assumed known, yielding $h_{1,t,\theta}^T$, second with base (\tilde{f}_θ^q) in the model with $\theta_1, \dots, \theta_{p'}$ known, yielding $\tilde{h}_{t,\theta}^T$.

We return now to the discussion of the optimal martingale estimating function (8) determined by the base $(f_\theta^1, \dots, f_\theta^r)$. Since for any $t = \Delta > 0$, $g_{\Delta,\theta}^{\text{opt}}$ is better than a $g_{\Delta,\theta}$ where $h_{\Delta,\theta} = h_\theta$ is determined as in Theorem 2, the flow $\mathcal{G}^{\text{opt}} = (g_{t,\theta}^{\text{opt}})$ should be small Δ -optimal if $r \geq d$ in case (i) or $r \geq \dim(d)$ in cases (ii) or (iii). What however does not follow from Theorem 2 is that \mathcal{G}^{opt} satisfies the conditions (11), (12) or (13). We shall now verify that this is the case (for $r = d$ and $\dim(d)$ respectively) and also argue that the lower bounds d and $\dim(d)$ for r cannot be improved upon.

Theorem 3 For the optimal flow $\mathcal{G}^{\text{opt}} = (g_{t,\theta}^{\text{opt}})$ of martingale estimating functions with base $(f_\theta^1, \dots, f_\theta^r)$ it holds that:

(i) If $r = d$ and the matrix $\partial_x f_\theta(x)$ is non-singular for μ_θ -a.a. x , then $g_{0,\theta}^{\text{opt}}(x, y) = \lim_{t \rightarrow 0} g_{t,\theta}^{\text{opt}}(x, y)$ and (11) holds,

$$\partial_y g_{0,\theta}^{\text{opt}}(x, x) = \dot{b}_\theta^T(x) C^{-1}(x).$$

(ii) If $r = \dim(d)$ and the matrix

$$\left(\partial_x f_\theta(x) \quad \partial_{xx}^2 f_\theta(x) \right) \in \mathbb{R}^{r \times (d+d^2)}$$

is of full rank $\dim(d)$ for μ_θ -a.a. x , then $g_{0,\theta}^{\text{opt}}(x, y) = \lim_{t \rightarrow 0} t g_{t,\theta}^{\text{opt}}(x, y)$ and (12) holds,

$$\partial_y g_{0,\theta}^{\text{opt}}(x, x) = 0, \quad \partial_{yy}^2 g_{0,\theta}^{\text{opt}}(x, x) = \dot{C}_\theta^T(x) (C_\theta(x)^{\otimes 2})^{-1}.$$

(iii) If $r = \dim(d)$ and the matrix

$$\left(\partial_x f_\theta(x) \quad \partial_{xx}^2 f_\theta(x) \right) \in \mathbb{R}^{r \times (d+d^2)}$$

is of full rank $\dim(d)$, then $g_{0,\theta}^{\text{opt}}(x, y) = \lim_{t \rightarrow 0} \begin{pmatrix} t g_{1,t,\theta}^{\text{opt}} \\ g_{2,t,\theta}^{\text{opt}} \end{pmatrix}$ and (13) holds,

$$\partial_y g_{0,\theta}^{\text{opt}}(x, x) = \begin{pmatrix} 0_{p' \times d} \\ \dot{b}_{2,\theta}^T(x) C_{\theta(x)}^{-1} \end{pmatrix}, \quad \partial_{yy}^2 g_{0,\theta}^{\text{opt}}(x, x) = \dot{C}_{1,\theta}^T(x) (C_\theta^{\otimes 2}(x))^{-1}. \quad (22)$$

If $r < d$ in (i) or $r < \dim(d)$ in (ii) or (iii), (11), resp. (12) or (13) will not be satisfied except possibly for a special choice of base $(f_\theta^1, \dots, f_\theta^r)$.

Proof. The main difficulty consists in finding $g_{0,\theta}^{\text{opt}}$ from (8) and (7). First note that for smooth functions ϕ, ψ , since

$$\begin{aligned} \pi_{t,\theta}(\phi\psi) &= \phi\psi + tA_\theta(\phi\psi) + o(t), \\ (\pi_{t,\theta}\phi)(\pi_{t,\theta}\psi) &= (\phi + tA_\theta\phi)(\psi + tA_\theta\psi) + o(t) \end{aligned}$$

(with $o(t)/t = o(t, x)/t \rightarrow 0$ for each x) and

$$A_\theta(\phi\psi) = (A_\theta\phi)\psi + \phi(A_\theta\psi) + (\partial_x\phi)C_\theta(\partial_x\psi)^T + o(t) \quad (23)$$

it follows that

$$\pi_{t,\theta}(\phi\psi) - (\pi_{t,\theta}\phi)(\pi_{t,\theta}\psi) = t(\partial_x\phi)C_\theta(\partial_x\psi)^T + o(t). \quad (24)$$

Now use (24) with $\phi = f^q$, $\psi = f^{q'}$, $1 \leq q, q' \leq r$ (with r arbitrary at the moment) to obtain

$$\pi_{t,\theta}(f_\theta f_\theta^T) - (\pi_{t,\theta}f_\theta)(\pi_{t,\theta}f_\theta)^T = t(\partial_x f_\theta)C_\theta(\partial_x f_\theta)^T + o(t) \quad (25)$$

and it is seen that the main term on the right evaluated at x is a non-singular $r \times r$ -matrix only if $r \leq d$ and $\partial_x f(x) \in \mathbb{R}^{r \times d}$ is of full rank r .

To find $g_{0,\theta}^{\text{opt}}$ we also need to approximate the factor $\partial_\theta \pi_{t,\theta} f_\theta - \pi_{t,\theta} \dot{f}$ from (7). Assuming that $\partial_\theta o(t) = o(t)$ and that \dot{f}_θ is smooth enough,

$$\begin{aligned} \partial_\theta \pi_{t,\theta} f_\theta &= \partial_\theta (f_\theta + tA_\theta f_\theta + o(t)) \\ &= \dot{f}_\theta + t \left(A_\theta \dot{f}_\theta + (\partial_x f_\theta) \dot{b}_\theta + \frac{1}{2} (\partial_{xx}^2 f_\theta) \dot{C}_\theta \right) + o(t), \\ \pi_{t,\theta} \dot{f}_\theta &= \dot{f}_\theta + tA_\theta \dot{f}_\theta + o(t) \end{aligned}$$

and thus

$$\partial_\theta \pi_{t,\theta} f_\theta - \pi_{t,\theta} \dot{f}_\theta = t \left((\partial_x f_\theta) \dot{b}_\theta + \frac{1}{2} (\partial_{xx}^2 f_\theta) \dot{C}_\theta \right) + o(t). \quad (26)$$

Case (i). Since $\dot{C}_\theta \equiv 0$, (26) reduces to

$$\partial_\theta \pi_{t,\theta} f_\theta - \pi_{t,\theta} \dot{f}_\theta = t(\partial_x f_\theta) \dot{b}_\theta + o(t),$$

and therefore, using (25) it follows that if $r \leq d$,

$$\begin{aligned} g_{0,\theta}^{\text{opt}}(x, y) &= \lim_{t \rightarrow 0} g_{t,\theta}^{\text{opt}}(x, y) \\ &= \dot{b}_\theta^T(x) (\partial_x f_\theta)^T(x) \left[\partial_x f_\theta(x) C(x) (\partial_x f_\theta)^T(x) \right]^{-1} (f_\theta(y) - f_\theta(x)) \end{aligned}$$

so that

$$\partial_y g_{0,\theta}^{\text{opt}}(x, x) = \dot{b}_\theta^T(x) (\partial_x f_\theta)^T(x) \left[\partial_x f_\theta(x) C(x) (\partial_x f_\theta)^T(x) \right]^{-1} \partial_x f_\theta(x). \quad (27)$$

For $r = d$ this reduces to $\dot{b}_\theta^T(x) C^{-1}(x)$ as wanted. For $r < d$ the $d \times d$ -matrix appearing as a factor to the right of $\dot{b}_\theta^T(x)$, has rank r , hence can never equal the non-singular $d \times d$ -matrix $C^{-1}(x)$. (However, it may still be possible to

obtain $\partial_y g_{0,\theta}^{\text{opt}}(x, x) = \dot{b}_\theta^T(x)C^{-1}(x)$: multiplying from the right by $C(x)$, it is seen that this holds for $\partial_y g_{0,\theta}^{\text{opt}}(x, x)$ given by (27) iff the $d \times d$ -matrix

$$(\partial_x f_\theta)^T(x) \left[\partial_x f_\theta(x)C(x) (\partial_x f_\theta)^T(x) \right]^{-1} \partial_x f_\theta(x)C(x)$$

acts as the identity on the subspace $L_b = \text{rowspan} \left(\dot{b}_\theta^T(x) \right)$ (multiplying by row vectors from the left). If $r < d$ this will be possible only if $\dim(L_b) < d$ (e.g. if $p < d$) and a special choice of f_θ determined by b_θ and C is required).

Case (ii). Assume that $r > d$. Then the main term on the right of (25) becomes singular and it is therefore necessary to expand further. But from the basic expansion

$$\pi_{t,\theta}\varphi = \varphi + tA_\theta\varphi + \frac{1}{2}t^2A_\theta^2\varphi + o(t^2),$$

using (23) repeatedly it eventually follows that

$$\pi_{t,\theta} (f_\theta f_\theta^T) - (\pi_{t,\theta} f_\theta) (\pi_{t,\theta} f_\theta)^T = t (\partial_x f_\theta) C_\theta (\partial_x f_\theta)^T + \frac{1}{2}t^2 Q + o(t^2) \quad (28)$$

with Q of the form

$$Q = (\partial_{xx}^2 f_\theta) C_\theta^{\otimes 2} (\partial_{xx}^2 f_\theta)^T + (\partial_x f_\theta) S + S^T (\partial_x f_\theta)^T \quad (29)$$

for some $S(x) \in \mathbb{R}^{d \times r}$. By Lemma 4 from the appendix therefore (with $A = (\partial_x f_\theta) C_\theta (\partial_x f_\theta)^T$, $B = \frac{1}{2}Q$)

$$\lim_{t \rightarrow 0} t^2 \left[\pi_{t,\theta} (f_\theta f_\theta^T) - (\pi_{t,\theta} f_\theta) (\pi_{t,\theta} f_\theta)^T \right]^{-1} = \mathcal{O}_2^T (\mathcal{O}_2 \frac{1}{2} Q \mathcal{O}_2^T)^{-1} \mathcal{O}_2 \quad (30)$$

where $\mathcal{O}(x) = \begin{pmatrix} \mathcal{O}_1(x) \\ \mathcal{O}_2(x) \end{pmatrix} \in \mathbb{R}^{r \times r}$ is orthogonal for each x , $\mathcal{O}_1(x)$ comprising the first d and $\mathcal{O}_2(x)$ the last $r - d$ rows of $\mathcal{O}(x)$, and satisfies

$$\mathcal{O}(x) (\partial_x f_\theta(x)) C_\theta(x) (\partial_x f_\theta)^T(x) \mathcal{O}^T(x) = \text{diag}(\lambda_1(x), \dots, \lambda_d(x), 0, \dots, 0) \quad (31)$$

with $\lambda_1(x), \dots, \lambda_d(x) > 0$ the non-zero eigenvalues for $(\partial_x f_\theta(x)) C_\theta(x) (\partial_x f_\theta)^T(x)$.

But from (31) follows that

$$\mathcal{O}_2(x) (\partial_x f_\theta(x)) C_\theta(x) (\partial_x f_\theta)^T(x) \mathcal{O}_2^T(x) = 0$$

or, since $C_\theta(x) \succ 0$, that

$$\mathcal{O}_2 \partial_x f_\theta = 0. \quad (32)$$

Combining (30) with (26) and using (32) it follows that

$$\begin{aligned} g_{0,\theta}^{\text{opt}}(x, y) &= \lim_{t \rightarrow 0} t g_{t,\theta}^{\text{opt}}(x, y) \\ &= \dot{C}_\theta^T (\partial_{xx}^2 f_\theta)^T \mathcal{O}_2^T \left[\mathcal{O}_2 (\partial_{xx}^2 f_\theta) C_\theta^{\otimes 2} (\partial_{xx}^2 f_\theta)^T \mathcal{O}_2^T \right]^{-1} \\ &\quad \times \mathcal{O}_2 (f_\theta(y) - f_\theta(x)) \end{aligned}$$

with all factors to the left of $f_\theta(y)$ evaluated at x . Using (32) it is clear that $\partial_y g_{0,\theta}^{\text{opt}}(x, x) = 0$ always and hence, to obtain (12), it remains to check whether (omitting the argument x with $\partial_{xx}^2 g_{0,\theta}^{\text{opt}}$ short for $\partial_{yy}^2 g_{0,\theta}^{\text{opt}}(x, x)$)

$$\partial_{xx}^2 g_{0,\theta}^{\text{opt}} = \dot{C}_\theta^T (\partial_{xx}^2 f_\theta)^T \mathcal{O}_2^T \left[\mathcal{O}_2 (\partial_{xx}^2 f_\theta) C_\theta^{\otimes 2} (\partial_{xx}^2 f_\theta)^T \mathcal{O}_2^T \right]^{-1} \mathcal{O}_2 \partial_{xx}^2 f_\theta \quad (33)$$

equals $\dot{C}_\theta^T (C_\theta^{\otimes 2})^{-1}$.

To achieve this we now assume that $r = \dim(d)$, so that $r - d = |J|$ and use the assumption from the theorem that $\partial_{xx}^2 f_\theta(x)$ has full rank $|J|$ for all x . Then $\Gamma := (\partial_{xx}^2 f_\theta) R$ also has rank $|J|$ and $\mathcal{O}_2 \Gamma \in \mathbb{R}^{J \times J}$ is non-singular and using that $\partial_{xx}^2 f_\theta = \partial_{xx}^2 f_\theta (R \tilde{R})$ (cf. (15)), (33) therefore gives

$$\begin{aligned} (\partial_{xx}^2 g_{0,\theta}^{\text{opt}}) R &= \dot{C}_\theta^T \tilde{R}^T \Gamma^T \mathcal{O}_2^T \left[\mathcal{O}_2 \Gamma \tilde{R} C_\theta^{\otimes 2} \tilde{R}^T \Gamma^T \mathcal{O}_2^T \right]^{-1} \mathcal{O}_2 \Gamma \\ &= \dot{C}_\theta^T \tilde{R}^T \left(\tilde{R} C_\theta^{\otimes 2} \tilde{R}^T \right)^{-1}. \end{aligned} \quad (34)$$

That $\partial_{xx}^2 g_{0,\theta}^{\text{opt}} = \dot{C}_\theta^T (C_\theta^{\otimes 2})^{-1}$ will follow from (cf. (15))

$$(\partial_{xx}^2 g_{0,\theta}^{\text{opt}}) R = \dot{C}_\theta^T (C_\theta^{\otimes 2})^{-1} R,$$

and that the right hand side here indeed equals that of (34) is verified multiplying by $\tilde{R} C_\theta^{\otimes 2} \tilde{R}^T$ from the right again appealing to (15).

Case (iii). Here we initially proceed as in case (ii), arriving at (cf. (28))

$$\begin{aligned} g_{t,\theta}^{\text{opt}}(x, y) &= \left(\dot{b}_\theta^T (\partial_x f_\theta)^T + \frac{1}{2} \dot{C}_\theta^T (\partial_{xx}^2 f_\theta)^T + o(1) \right) \\ &\quad \times \left((\partial_x f_\theta) C_\theta (\partial_x f_\theta)^T + \frac{1}{2} t Q + o(t) \right)^{-1} (f_\theta(y) - f_\theta(x) + o(1)) \end{aligned} \quad (35)$$

with Q as in (29).

Considering first the last $p - p'$ components of $g_{t,\theta}^{\text{opt}}$, since by assumption $\dot{C}_{2,\theta} \equiv 0$, it follows from Lemma 4 that

$$g_{2,t,\theta}^{\text{opt}}(x, y) = b_{2,\theta}^T (\partial_x f_\theta)^T \left(\frac{1}{t} \mathcal{O}_2^T (\mathcal{O}_2 \frac{1}{2} Q \mathcal{O}_2^T)^{-1} \mathcal{O}_2 + N \right) \times (f_\theta(y) - f_\theta(x)) + o(1),$$

which because of (32) in the limit reduces to

$$g_{2,0,\theta}^{\text{opt}}(x, y) = \lim_{t \rightarrow 0} g_{2,t,\theta}^{\text{opt}}(x, y) = b_{2,\theta}^T (\partial_x f_\theta)^T N (f_\theta(y) - f_\theta(x))$$

with N of the form

$$N = \mathcal{O}_1^T \left(\mathcal{O}_1 (\partial_x f_\theta) C_\theta (\partial_x f_\theta)^T \mathcal{O}_1^T \right)^{-1} \mathcal{O}_1 + \mathcal{O}_2^T \tilde{S} + \tilde{S}^T \mathcal{O}_2.$$

But then, again using (32) and since $\mathcal{O}_1 (\partial_x f_\theta) \in \mathbb{R}^{d \times d}$ is non-singular

$$\begin{aligned} \partial_y g_{2,0,\theta}^{\text{opt}}(x, x) &= b_{2,\theta}^T (\partial_x f_\theta)^T \mathcal{O}_1^T \left(\mathcal{O}_1 (\partial_x f_\theta) C_\theta (\partial_x f_\theta)^T \mathcal{O}_1^T \right)^{-1} \mathcal{O}_1 (\partial_x f_\theta) \\ &= b_{2,\theta}^T C_\theta^{-1} \end{aligned}$$

as wanted in the first part of (22).

As for the first p' components of $g_{t,\theta}^{\text{opt}}$, obtain from (35) that

$$tg_{1,t,\theta}^{\text{opt}}(x, y) = \left(b_{1,\theta}^T (\partial_x f_\theta)^T + \frac{1}{2} \dot{C}_{1,\theta}^T (\partial_{xx}^2 f_\theta)^T \right) \mathcal{O}_2^T (\mathcal{O}_2 \frac{1}{2} Q \mathcal{O}_2^T)^{-1} \mathcal{O}_2 \times (f_\theta(y) - f_\theta(x)) + o(1),$$

whence

$$\begin{aligned} g_{1,0,\theta}^{\text{opt}}(x, y) &= \lim_{t \rightarrow 0} tg_{1,t,\theta}^{\text{opt}}(x, y) \\ &= \frac{1}{2} \dot{C}_{1,\theta}^T (\partial_{xx}^2 f_\theta)^T \mathcal{O}_2^T \left(\mathcal{O}_2 \frac{1}{2} (\partial_{xx}^2 f_\theta) C^{\otimes 2} (\partial_{xx}^2 f_\theta)^T \mathcal{O}_2^T \right)^{-1} \\ &\quad \times \mathcal{O}_2 (f_\theta(y) - f_\theta(x)) \end{aligned}$$

once again using (32). But then (32) also gives

$$\partial_y g_{1,0,\theta}^{\text{opt}}(x, x) = 0.$$

and arguing exactly as in the last part of case (ii), one finally finds that

$$\partial_{yy}^2 g_{1,0,\theta}^{\text{opt}}(x, x) = \dot{C}_{1,\theta}^T (C_\theta^{\otimes 2})^{-1},$$

and we have completed the proof of (22). ■

We have not shown that (11), respectively (12), (13) are satisfied for the optimal martingale estimating function when $r > d$, resp. $r > \dim(d)$. For case (i) with $r > d$ one may copy the argument involving $g_{2,\theta}^{\text{opt}}$ given in case (iii) above. For cases (ii) and (iii), if $r > \dim(d)$ a further expansion of (28) together with a refinement of the lemma is required, since e.g. the columns of $A = (\partial_x f_\theta) C_\theta (\partial_x f_\theta)^T$ and $B = \frac{1}{2}Q$ cannot span a subspace of dimension r . We believe however that the relevant of (12) or (13) is still valid for the optimal martingale estimating function, even if $r > \dim(d)$.

Remark 3 For $d = 1$, Bibby and Sørensen [1] studied martingale estimating functions with the one-dimensional ($r = 1$) base $f(x) = x$, and apart from deriving the optimal estimating function G_n^* ((2.15) in [1]), also suggested the use of the approximately optimal \tilde{G}_n ((2.14) in [1]). As they point out, the weights for \tilde{G}_n are arrived at by replacing the true transition probabilities as they appear in (7) by the Gaussian approximations corresponding to the Euler scheme, i.e. the conditional distribution of X_t given $X_0 = x$ is approximated by the normal distribution $n_{t,\theta}(x, \cdot)$ with mean $x + tb_\theta(x)$ and variance $tC_\theta(x)$. We shall now discuss how the use of this approximation may still lead to estimating functions that are small Δ -optimal. We shall only consider one-dimensional diffusions but allow p and r to be arbitrary.

In terms of $n_{t,\theta}(x, \cdot)$, the approximation to $\pi_{t,\theta}\phi$ becomes

$$\begin{aligned} \pi_{t,\theta}\phi(x) &\approx \int n_{t,\theta}(x, dy) \phi(y) \\ &= \phi(x) + t(b_\theta(x)\phi'(x) + \frac{1}{2}C_\theta(x)\phi''(x)) \\ &\quad + \frac{1}{2}t^2(b_\theta^2(x)\phi''(x) + b_\theta(x)C_\theta(x)\phi'''(x) + \frac{1}{4}C_\theta^2(x)\phi''''(x)) + o(t^2) \end{aligned}$$

resulting in the approximations (omitting the argument x),

$$\pi_{t,\theta}(\phi\psi) - (\pi_{t,\theta}\phi)(\pi_{t,\theta}\psi) \approx n_{t,\theta}(\phi, \psi)$$

where

$$\begin{aligned} n_{t,\theta}(\phi, \psi) &= tC_\theta\phi'\psi' \\ &\quad + \frac{1}{2}t^2[C_\theta^2\phi''\psi'' + 2b_\theta C_\theta(\phi''\psi' + \phi'\psi'')] \\ &\quad + C_\theta^2(\phi'''\psi' + \phi'\psi''') + o(t^2), \end{aligned} \tag{36}$$

and

$$\partial_\theta (\pi_{t,\theta} \phi) \approx t \left(\dot{b}_\theta \phi' + \frac{1}{2} \dot{C}_\theta \phi'' \right) + o(t).$$

Given a base $(f^q)_{1 \leq q \leq r}$ (for convenience assumed not to depend on θ), now consider the martingale estimating function (cf. (7) and (8))

$$g_{t,\theta}^{\text{norm}}(x, y) = (h_{t,\theta}^{\text{norm}})^T(x) (f(y) - \pi_{t,\theta} f(x)) \quad (37)$$

where

$$h_{t,\theta}^{\text{norm}} = (N_{t,\theta} f)^{-1} t \left(\partial_x f \dot{b}_\theta + \frac{1}{2} \partial_{xx}^2 f \dot{C}_\theta \right)$$

writing $N_{t,\theta} f$ for the $r \times r$ -matrix with (q, q') 'th element $n_{t,\theta}(f^q, f^{q'})$. From the approximations above it follows that

$$N_{t,\theta} f = t (\partial_x f) C_\theta (\partial_x f)^T + \frac{1}{2} t^2 \tilde{Q} + o(t^2)$$

with \tilde{Q} of the same form as Q in (29), and this is enough for the proof of Theorem 3 to carry over and yield the following result: the estimating function (37) is small Δ -optimal (in the sense of (11), (12) or (13)) in case (i) if $r = 1$ and in cases (ii) or (iii) if $r = \dim(1) = 2$. In case (i) with $r = 1$, instead of using $N_{t,\theta} f = n_{t,\theta}(f, f)$ one may use

$$N_{t,\theta} f = t C_\theta (\partial_x f)^2$$

corresponding to the main term in (36). In cases (ii) and (iii) it is essential to use (36) as it stands and that $r \geq 2$.

3 Examples

We shall illustrate the results of the previous sections through two examples.

3.1 A generalized Cox-Ingersoll-Ross process

Consider the one-dimensional ($d = 1$) SDE

$$dX_t = (aX_t^{2\gamma-1} + bX_t) dt + \sigma X_t^\gamma dB_t \quad (38)$$

where $a, b \in \mathbb{R}$, $\gamma \neq 1$ and $\sigma > 0$. For $\gamma = \frac{1}{2}$ this is the SDE for the Cox-Ingersoll-Ross process (CIR-process, see (39) below). The generalization

(38) is arrived at by considering all powers \tilde{X}^ρ of a CIR with $\rho \neq 0$, more precisely if X solves (38), then the associated CIR-process is $\tilde{X} = X^{2-2\gamma}$ solving

$$d\tilde{X}_t = \left(\tilde{a} + \tilde{b}X_t \right) dt + \tilde{\sigma} \sqrt{\tilde{X}_t} dB_t \quad (39)$$

where

$$\tilde{b} = (2 - 2\gamma)b, \quad \tilde{\sigma}^2 = (2 - 2\gamma)^2 \sigma^2, \quad \tilde{a} - \frac{1}{2}\tilde{\sigma}^2 = (2 - 2\gamma) \left(a - \frac{1}{2}\sigma^2 \right), \quad (40)$$

(which also explains why $\gamma = 1$ is not allowed in (38)).

Because of the connection to the CIR-process, the model described by (38) is much simpler to handle than the more standard CKLS-model,

$$dX_t = (a + bX_t) dt + \sigma X_t^\gamma dB_t,$$

in particular, for (38) it is easy to find martingale estimating functions of the type considered in the preceding sections.

In (38) the parameter space has dimension $p = 4$. We shall want X to be strictly positive and ergodic, which happens iff the associated CIR-process \tilde{X} is strictly positive and ergodic, i.e. $\tilde{b} < 0$ and $2\tilde{a} \geq \tilde{\sigma}^2$, or equivalently, either $\gamma < 1$, $b < 0$, $2a \geq \sigma^2$ or $\gamma > 1$, $b > 0$, $2a \leq \sigma^2$. As *open* parameter set we shall therefore use

$$\Theta = \left\{ (a, b, \gamma, \sigma^2) : \sigma^2 > 0 \text{ and } \begin{array}{l} \text{either } \gamma < 1, b < 0, 2a > \sigma^2 \\ \text{or } \gamma > 1, b > 0, 2a < \sigma^2 \end{array} \right\}.$$

Note that if $\theta = (a, b, \gamma, \sigma^2) \in \Theta$ and $\rho \neq 0$, then X^ρ solves (38) with parameters $\theta^* = (a^*, b^*, \gamma^*, \sigma^{*2})$ given by

$$b^* = \rho b, \quad \sigma^{*2} = \rho^2 \sigma^2, \quad 2 - 2\gamma^* = \frac{1}{\rho} (2 - 2\gamma), \quad a^* - \frac{1}{2}\sigma^{*2} = \rho \left(a - \frac{1}{2}\sigma^2 \right).$$

In particular, taking $\rho < 0$ corresponds to a switch from $\gamma < 1$ to $\gamma^* > 1$ (or from $\gamma > 1$ to $\gamma^* < 1$).

Since the invariant distribution for \tilde{X} is a Γ -distribution, the invariant distribution for X is that of a Γ -distributed random variable raised to the power $(2 - 2\gamma)^{-1}$. The density is

$$\mu_\theta(x) = \frac{|2 - 2\gamma|}{\Gamma\left(\frac{2\tilde{a}}{\tilde{\sigma}^2}\right) \left(\frac{(2\gamma-2)\sigma^2}{2b}\right)^{2\tilde{a}/\tilde{\sigma}^2}} x^{\frac{2a}{\sigma^2}-2\gamma} \exp\left(-\frac{2b}{(2\gamma-2)\sigma^2} x^{2-2\gamma}\right) \quad (41)$$

for $x > 0$, where (cf (40))

$$\frac{2\tilde{a}}{\tilde{\sigma}^2} = \frac{2a}{(2-2\gamma)\sigma^2} + \frac{1-2\gamma}{2-2\gamma}.$$

(For $\gamma = \frac{1}{2}$ the familiar invariant Γ -density for the CIR-process is obtained).

Because a Γ -distribution has finite moments of all orders $m \in \mathbb{N}$ we have $E_\theta^\mu X_0^{(2\gamma-2)m} < \infty$ for all $m \in \mathbb{N}$, and since

$$E_\theta^\mu X_t^{(2\gamma-2)m} = \int_0^\infty dx \mu_\theta(x) \pi_{t,\theta} x^{(2\gamma-2)m}$$

(where $\pi_{t,\theta} x^\beta$ is short for $\pi_{t,\theta} f(x)$ for $f(y) = y^\beta$), also

$$\pi_{t,\theta} x^{(2\gamma-2)m} < \infty$$

for all $t > 0$, $m \in \mathbb{N}$ and (Lebesgue almost all) $x > 0$.

The conditional moments for a CIR-process are known and in any case easy to find using polynomial martingales: for $m \in \mathbb{N}$, let $\tilde{\xi}_m$ be the m 'th moment in the invariant distribution for \tilde{X} ,

$$\tilde{\xi}_m = \left(-\frac{\tilde{\sigma}^2}{2\tilde{b}} \right)^m \frac{\Gamma\left(\frac{2\tilde{a}}{\tilde{\sigma}^2} + m\right)}{\Gamma\left(\frac{2\tilde{a}}{\tilde{\sigma}^2}\right)}$$

and verify, for instance using induction on m and Itô's formula, that $M^{(m)}$ is a mean-zero martingale (see the note below) under each P_θ^x , where

$$M_t^{(m)} = e^{-\tilde{b}mt} \sum_{i=1}^m \beta_i^{(m)} \left(X_t^{(2-2\gamma)i} - \tilde{\xi}_i \right) \quad (42)$$

with

$$\beta_i^{(m)} = \frac{1}{\tilde{\xi}_i} (-1)^{i-1} \binom{m}{i}.$$

(Equivalently, for each m , the polynomial $\sum_i \beta_i^{(m)} (x^i - \tilde{\xi}_i)$ of degree m is an eigenfunction for the generator for the CIR-process (39) corresponding to the eigenvalue $\tilde{b}m$, see Kessler and Sørensen [5] for estimating functions built from eigenfunctions, and their Example 2.1 for the CIR-process).

Note. Because all conditional moments for the ergodic CIR-process are finite, one verifies directly that the local martingale $M^{(m)}$ satisfies $E_\theta^x [M^{(m)}]_t <$

∞ for all x and t , in particular $M^{(m)}$ is a true martingale under P_θ^x (L^2 -bounded on $[0, t]$ for all t).

Turning now to the problem of estimating θ from discrete observations of X , it is clear that the model (38) belongs to case (iii) with $p = 4$, $p' = 2$, so we shall apply Theorems 2 and 3 for that case with $r = 1 + \dim(1) = 2$. In view of the above, a simple candidate for the base (f^1, f^2) is

$$f^1(x) = x^{2-2\gamma}, \quad f^2(x) = x^{4-4\gamma}, \quad (43)$$

which trivially satisfies Assumption A. Note that f^1, f^2 both depend on θ , cf. the comment immediately preceding Proposition 1.

Using Theorem 2 with $*$ in (18) equal to 0, and listing the parameters in the order γ, σ^2, a, b one finds that

$$\dot{b}_{2,\theta}^T(x) = \begin{pmatrix} x^{2\gamma-1} \\ x \end{pmatrix}, \quad \dot{C}_{1,\theta}^T(x) = \begin{pmatrix} 2\sigma^2 x^{2\gamma} \log x \\ x^{2\gamma} \end{pmatrix}$$

and eventually arrives at the estimating function

$$g_{t,\theta}(x, y) = \begin{pmatrix} -2 \log x & x^{2\gamma-2} \log x \\ -2 & x^{2\gamma-2} \\ 2(3-4\gamma)x^{2\gamma-2} & -(1-2\gamma)x^{4\gamma-4} \\ 2(3-4\gamma) & -(1-2\gamma)x^{2\gamma-2} \end{pmatrix} \begin{pmatrix} y^{2-2\gamma} - \pi_{t,\theta} x^{2-2\gamma} \\ y^{4-4\gamma} - \pi_{t,\theta} x^{4-4\gamma} \end{pmatrix}, \quad (44)$$

which requires the use of (42) for $m = 1, 2$ in order to find the conditional expectations.

That $g_{t,\theta}$ given by (44) indeed satisfies the conditions (13) for small Δ -optimality is most easily verified directly. Note that the asserted (Theorem 2) linear independence between the columns in h_θ , i.e. the functions comprising the rows in the 4×2 -matrix in (44) where rows 2 and 4 are similar, holds precisely because $\gamma \neq 1$. Note that this similarity also implies that the 4×2 -matrix in (44) may be replaced the simpler

$$\begin{pmatrix} -2 \log x & x^{2\gamma-2} \log x \\ 1 & 0 \\ 2(3-4\gamma)x^{2\gamma-2} & -(1-2\gamma)x^{4\gamma-4} \\ 0 & x^{2\gamma-2} \end{pmatrix}$$

without affecting the solutions to the estimating equations.

For the flow $(g_{t,\theta})$ given by (44) one still needs to check the integrability assumptions from Jacobsen [2] Theorem 7.5, and the conditions on estimating flows made prior to that theorem. In the case at hand the conditions in particular amount to requiring that $E_\theta^\mu |g_{t,\theta}^k(X_0, X_t)|^K < \infty$ for all components k and moderate values of $K \in \mathbb{N}$. The problem is the appearance of powers $x^{2\gamma-2}$ and $x^{4\gamma-4}$ in the expression for $g_{t,\theta}$, which translates into negative powers \tilde{X}_0^{-1} and \tilde{X}_0^{-2} of the CIR-process \tilde{X} , and of course e.g. $E_\theta^\mu \tilde{X}_0^{-K} = E_\theta^\mu X_0^{(2\gamma-2)K} < \infty$ iff $\frac{2\tilde{a}}{\tilde{\sigma}^2} > K$. Thus some care should be taken before applying (44), at least it must be assumed that $\frac{2\tilde{a}}{\tilde{\sigma}^2}$ be suitably large.

To find the optimal martingale estimating function with base f^1, f^2 given by (43), one needs (42) also for $m = 3, 4$ and conditional moments involving logarithms, see (7) in which the term $\pi_{t,\theta} \dot{f}_\theta$ appears. The latter moments are easy to find in terms of the conditional expectation $E_\theta(\log \tilde{X}_t | \tilde{X}_0 = \tilde{x})$ for the CIR-process starting at an arbitrary level \tilde{x} , but the explicit form for this is unpleasant to work with.

Whether one uses the small Δ -optimal flow (44) or the optimal flow, since $d = 1$ a slight improvement in efficiency may be gained by symmetrizing, using e.g. $\frac{1}{2}(g_{t,\theta}(x, y) + g_{t,\theta}(y, x))$ instead of (44), cf. Jacobsen [2], Proposition 6.1, and the discussion there about time reversal.

3.2 The finite-dimensional Gaussian diffusions

We consider now the d -dimensional diffusion

$$dX_t = (A + BX_t) dt + D dW_t \quad (45)$$

where the unknown parameters are $A \in \mathbb{R}^{d \times 1}$, $B \in \mathbb{R}^{d \times d}$ and $C := DD^T \in \mathbb{R}^{d \times d}$, with the symmetric matrix C assumed strictly positive definite. (In this subsection the symbol B is used to denote the matrix of linear drift parameters and the driving d -dimensional Brownian motion is denoted W instead of B). Thus

$$p' = |J|, \quad p = |J| + d + d^2.$$

The diffusion (45) has Gaussian transitions (for the expectation and second order moments, see (46) and (47) below) and is ergodic iff $\text{spec}(B) \subset \{\lambda \in \mathbb{C} : \text{Re}(\lambda) < 0\}$.

As base (of dimension $\dim(d)$) for the martingale estimating functions we shall use (f^q) for $q \in \{1, \dots, d\} \cup J$, where

$$f^i(x) = x_i \quad (1 \leq i \leq d), \quad f^{i'j'}(x) = x_{i'}x_{j'} \quad ((i', j') \in J),$$

writing $x = (x_1, \dots, x_d)$ for a generic point in \mathbb{R}^d . Clearly (f^q) satisfies the conditions from Assumption A and also, as a little work shows, the conditions on $\partial_x f$, $\partial_{xx}^2 f$ from Theorem 2. To proceed we need the conditional moments $\pi_{t,\theta} f^q$, conveniently collected in the vector $\pi_{t,\theta} x = (\pi_{t,\theta} x_i)_{1 \leq i \leq d}$ and the matrix $\pi_{t,\theta} x x^T$ and known to be given by the expressions

$$\pi_{t,\theta} x = (e^{tB} - I_d) B^{-1} A + e^{tB} x \quad (46)$$

$$\pi_{t,\theta} x x^T = (\pi_{t,\theta} x) (\pi_{t,\theta} x)^T + \int_0^t e^{sB} C e^{sB^T} ds. \quad (47)$$

(As in the previous example, notation like $\pi_{t,\theta} x x^T$ is short for $\pi_{t,\theta} f(x)$, where $f(y) = y y^T$.)

Invoking Theorem 2 with $*$ in (18) equal to 0, one eventually arrives at the following small Δ -optimal estimating function $g_{t,\theta}$ with $g_{1,t,\theta} = (g_{t,\theta}^{i'j'})_{(i',j') \in J}$ and $g_{2,t,\theta}$ split into the vector-valued component $g_{2,t,\theta}^A = (g_{t,\theta}^i)_{1 \leq i \leq d}$ and the matrix-valued component $g_{2,t,\theta}^B = (g_{t,\theta}^{ij})_{1 \leq i,j \leq d}$ and $g_{1,t,\theta}$ and $g_{2,t,\theta}^A$ and $g_{2,t,\theta}^B$ given by

$$\begin{aligned} g_{1,t,\theta}^{i'j'}(x, y) &= \left(C^{-1} \left[-x (y - \pi_{t,\theta} x)^T - (y - \pi_{t,\theta} x) x^T \right. \right. \\ &\quad \left. \left. + y y^T - \pi_{t,\theta} (x x^T) \right] C^{-1} \right)_{i'j'}, \\ g_{2,t,\theta}^A(x, y) &= C^{-1} (y - \pi_{t,\theta} x), \\ g_{2,t,\theta}^B(x, y) &= C^{-1} (y - \pi_{t,\theta} x) x^T. \end{aligned}$$

For the calculations one uses that

$$\dot{C}_{1,\theta}^T(x) \in \mathbb{R}^{|J| \times d^2}, \quad \left(\dot{C}_{1,\theta}^T(x) \right)_{i'j',ij} = \begin{cases} \delta_{i'i} \delta_{j'j} & \text{if } i \leq j, \\ \delta_{i'j} \delta_{j'i} & \text{if } i > j, \end{cases}$$

$$\left(\dot{b}_{2,\theta}^A \right)^T(x) = I_d \in \mathbb{R}^{d \times d},$$

$$\left(\dot{b}_{2,\theta}^B \right)^T(x) = I_d \otimes x \in \mathbb{R}^{d^2 \times d}.$$

Also note that

$$(\partial_x f(x) \quad \partial_{xx}^2 f(x) R)^{-1} = \begin{pmatrix} I_d & 0_{d \times J} \\ P(x) & D \end{pmatrix}$$

where $D = \text{diag}(d_{i'j'}) \in \mathbb{R}^{J \times J}$ with

$$d_{i'j'} = \begin{cases} 1 & \text{if } i' < j' \\ \frac{1}{2} & \text{if } i' = j', \end{cases}$$

and $P(x) \in \mathbb{R}^{J \times d}$ with

$$P_{i'j',j}(x) = -d_{i'j'} (\delta_{i'j} x_{j'} + \delta_{j'j} x_{i'}).$$

As in the previous example, the simplest way to verify the small Δ -optimality is to verify directly from these expressions that the conditions (13) are satisfied.

The resulting estimating equations are not affected by multiplication from the left and/or right by C , and it is now an easy task to write down the estimators of the parameter functions

$$\mathfrak{A} := (e^{\Delta B} - I_d) B^{-1} A, \quad e^{\Delta B}, \quad \mathfrak{C} := \int_0^\Delta e^{sB} C e^{sB^T} ds$$

based on the observations $X_0, X_\Delta, \dots, X_{n\Delta}$: defining

$$\bar{X}_* := \frac{1}{n} \sum_{i=1}^n X_{(i-1)\Delta}, \quad \bar{X}^* := \frac{1}{n} \sum_{i=1}^n X_{i\Delta},$$

using (46) and (47) one sees that the estimating equations obtained from g_1, g_2^A, g_2^B are equivalent to the equations

$$\sum_{i=1}^n (X_{i\Delta} - \mathfrak{A} - e^{\Delta B} X_{(i-1)\Delta}) = 0, \tag{48}$$

$$\sum_{i=1}^n (X_{i\Delta} - \mathfrak{A} - e^{\Delta B} X_{(i-1)\Delta}) X_{(i-1)\Delta}^T = 0, \tag{49}$$

$$\sum_{i=1}^n \left(X_{i\Delta} X_{i\Delta}^T - (\mathfrak{A} + e^{\Delta B} X_{(i-1)\Delta}) (\mathfrak{A} + e^{\Delta B} X_{(i-1)\Delta})^T - \mathfrak{C} \right) = 0, \tag{50}$$

and hence

$$\hat{\mathfrak{A}} = \bar{X}^* - e^{\Delta \hat{B}} \bar{X}_*, \quad (51)$$

$$e^{\Delta \hat{B}} = \left(\sum_{i=1}^n (X_{i\Delta} - \bar{X}^*) X_{(i-1)\Delta}^T \right) \left(\sum_{i=1}^n (X_{(i-1)\Delta} - \bar{X}_*) (X_{(i-1)\Delta} - \bar{X}_*)^T \right)^{-1} \quad (52)$$

$$\hat{\mathfrak{C}} = \frac{1}{n} \sum_{i=1}^n (X_{i\Delta} X_{i\Delta}^T - Z_i Z_i^T) \quad (53)$$

where in the last line,

$$Z_i := \hat{\mathfrak{A}} + e^{\Delta \hat{B}} X_{(i-1)\Delta}.$$

Note that the equation (53) may be written

$$\hat{\mathfrak{C}} = \frac{1}{n} \sum_{i=1}^n \left(X_{i\Delta} - \hat{\mathfrak{A}} - e^{\Delta \hat{B}} X_{(i-1)\Delta} \right) \left(X_{i\Delta} - \hat{\mathfrak{A}} - e^{\Delta \hat{B}} X_{(i-1)\Delta} \right)^T \quad (54)$$

as is seen using that from (48) and (49) it follows that

$$\sum_{i=1}^n \left(X_{i\Delta} - \hat{\mathfrak{A}} - e^{\Delta \hat{B}} X_{(i-1)\Delta} \right) Z_i^T = 0.$$

The likelihood function for observing $X_0, X_\Delta, \dots, X_{n\Delta}$ conditionally on X_0 is

$$\prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\mathfrak{C}|} \exp \left(-\frac{1}{2} (X_{i\Delta} - \xi_i)^T \mathfrak{C}^{-1} (X_{i\Delta} - \xi_i) \right)$$

where

$$\xi_i = \pi_{\Delta, \theta} (X_{(i-1)\Delta}) = \mathfrak{A} + e^{\Delta B} X_{(i-1)\Delta}.$$

Maximizing this over \mathfrak{A} , $e^{\Delta B}$ and \mathfrak{C} varying *freely* in $\mathbb{R}^{d \times 1}$, $\mathbb{R}^{d \times d}$ and the space of symmetric positive definite $d \times d$ - matrices yields the estimators $\hat{\mathfrak{A}}$, $e^{\Delta \hat{B}}$ and $\hat{\mathfrak{C}}$ from (51), (52) and (53). In a forthcoming paper by Mathieu Kessler and Anders Rahbek [4] (for the model with $A = 0$), the authors study this maximum-likelihood estimator and also tackle the non-trivial problem of converting the expressions (51), (52) and (53) into estimators for A , B and C : it may of course happen with probability > 0 that the right-hand side of (52) is not the exponential of any square matrix, and even if it is, \hat{B} may not satisfy the basic condition for ergodicity that $\text{Re}(\lambda) < 0$ for all $\lambda \in \text{spec}(\hat{B})$.

As pointed out by Kessler and Rahbek, a particularly unpleasant problem is that the \widehat{B} solving (52) may not even be unique, due to certain periodicities. If \widehat{B} is found of course

$$\widehat{A} = \widehat{B} \left(e^{\Delta \widehat{B}} - I_d \right)^{-1} \widehat{\mathfrak{A}}$$

(with $e^{\Delta \widehat{B}} - I_d$ non-singular iff 0 is not an eigenvalue for \widehat{B} , hence non-singular almost surely if only n is large enough). Finally \widehat{C} should be found from (53), but again here, Kessler and Rahbek show that there need not be a unique solution.

Remark 4 *If the observations X_{t_i} are not equidistant as assumed above (where $t_i = i\Delta$), maximum-likelihood estimation if at all possible will certainly be very difficult. Also, there are no analogues of $\widehat{\mathfrak{A}}$, $e^{\Delta \widehat{B}}$ and $\widehat{\mathfrak{C}}$ and in order to solve the estimating equations which are obtained from (48), (49) and (50), one must rely entirely on numerics and proceed directly to find \widehat{A} , \widehat{B} and \widehat{C} . A possible method is to reparametrize B through its spectrum and look for B only with d distinct eigenvalues (it may be argued that with probability one, \widehat{B} will have this property) that are either real or come in complex conjugate pairs. The parametrization of B could then be in terms of these eigenvalues and a concrete choice of corresponding complex eigenvectors of unit length. At least, if the t_i are non-lattice, the uniqueness problems emphasized by Kessler and Rahbek should disappear.*

A Appendix

The following result was used in the proof of Theorem 3:

Lemma 4 *Let $A, B \in \mathbb{R}^{m \times m}$ be symmetric and positive semidefinite matrices such that $1 \leq \text{rank}(A) = m' < m$ and such that the columns (or rows) of A and B jointly span all of \mathbb{R}^m . Let further $\mathcal{O} = \begin{pmatrix} \mathcal{O}_1 \\ \mathcal{O}_2 \end{pmatrix}$ be an orthogonal $m \times m$ -matrix with \mathcal{O}_1 comprising the first m' , \mathcal{O}_2 the last $m - m'$ rows of \mathcal{O} such that*

$$\mathcal{O}A\mathcal{O}^T = \text{diag}(\lambda_1, \dots, \lambda_{m'}, 0, \dots, 0),$$

$\lambda_1, \dots, \lambda_{m'} > 0$ denoting the non-zero eigenvalues for A . Then as $t \rightarrow 0$,

$$(A + tB)^{-1} = \frac{1}{t} \mathcal{O}_2^T (\mathcal{O}_2 B \mathcal{O}_2^T)^{-1} \mathcal{O}_2 + N + O(t),$$

where N is of the form

$$\mathcal{O}_1^T (\mathcal{O}_1 A \mathcal{O}_1^T)^{-1} \mathcal{O}_1 + \mathcal{O}_2^T S + S^T \mathcal{O}_2$$

for some $(m - m') \times m$ -matrix S .

Proof. Assume first that $A = \text{diag}(\lambda_1, \dots, \lambda_{m'}, 0, \dots, 0)$ (with all $\lambda_\ell > 0$) and write

$$B = \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$$

with e.g. B_{22} the lower right $(m - m') \times (m - m')$ -submatrix of B . Then

$$|A + tB| = t^{m-m'} \left(\prod_{\ell=1}^{m'} \lambda_\ell \right) |B_{22}| + O(t^{m-m'+1}),$$

Also for the subdeterminants obtained by deleting the ℓ 'th row and ℓ' 'th column,

$$|A + tB|_{\ell\ell'} = \begin{cases} O(t^{m-m'-1}) & \text{if } \ell, \ell' > m' \\ O(t^{m-m'}) & \text{otherwise.} \end{cases}$$

It follows from this that $(A + tB)^{-1}$ is of the form

$$\frac{1}{t} \begin{pmatrix} 0 & 0 \\ 0 & M \end{pmatrix} + N + O(t)$$

and it is then easy to see that, writing

$$D = \text{diag}(\lambda_1, \dots, \lambda_{m'}) \in \mathbb{R}^{m' \times m'}$$

one has

$$(A + tB)^{-1} = \frac{1}{t} \begin{pmatrix} 0 & 0 \\ 0 & B_{22}^{-1} \end{pmatrix} + \begin{pmatrix} D^{-1} & -D^{-1} B_{12} B_{22}^{-1} \\ -B_{22}^{-1} B_{21} D^{-1} & 0 \end{pmatrix} + O(t). \quad (55)$$

For the general case, just use that

$$(A + tB)^{-1} = \mathcal{O}^T (\mathcal{O} A \mathcal{O}^T + t \mathcal{O} B \mathcal{O}^T)^{-1} \mathcal{O}$$

with $(\mathcal{O} A \mathcal{O}^T + t \mathcal{O} B \mathcal{O}^T)^{-1}$ of the form (55) and $D = \mathcal{O}_1 A \mathcal{O}_1^T$. ■

References

- [1] Bibby, B.M., Sørensen, M. (1995). Martingale estimating functions for discretely observed diffusion processes. *Bernoulli* **1**, 17–39.
- [2] Jacobsen, M. (1998). Discretely observed diffusions: classes of estimating functions and small Δ –optimality. Preprint 11, Department of Theoretical Statistics, University of Copenhagen. (To appear in *Scand. J. Statist.*).
- [3] Kessler, M. (2000). Simple and explicit estimating functions for a discretely observed diffusion process. *Scand. J. Statist.* **27**, 65–82.
- [4] Kessler, M., Rahbek, A. (2000). Inference for discretely observed multivariate homogeneous Gaussian diffusions (manuscript).
- [5] Kessler, M., Sørensen, M. (1999). Estimating equations based on eigenfunctions for a discretely observed diffusion process. *Bernoulli* **5**, 299–314.
- [6] Sørensen, M. (1998). On asymptotics of estimating functions. Preprint 6, Department of Theoretical Statistics, University of Copenhagen.