

# On estimating a dynamic function of a stochastic system with averaging

Robert Liptser ([liptser@eng.tau.ac.il](mailto:liptser@eng.tau.ac.il))

*Dept. Electrical Engineering-Systems, Tel Aviv University, 69978 Tel Aviv, Israel*

Vladimir Spokoiny ([spokoiny@wias-berlin.de](mailto:spokoiny@wias-berlin.de))

*Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstr. 39, 10117 Berlin, Germany*

**Abstract.** We consider a two-scaled diffusion system, when drift and diffusion parameters of the “slow” component are contaminated by the “fast” unobserved component. The goal is to estimate the dynamic function which is defined by averaging the drift coefficient of the “slow” component w.r.t. the stationary distribution of the “fast” one. We apply a locally linear smoother with a data-driven bandwidth choice. The procedure is fully adaptive and nearly optimal up to a  $\log \log$  factor.

**Keywords:** fast and slow components, drift and diffusion coefficients, averaging principle, nonparametric estimation, bandwidth selection

**AMS Subclass 1991: Primary:** 62G05; Secondary 62M99

## 1. Introduction

In this paper, we propose a procedure for adaptive estimation of “averaged” characteristics of a two scaled diffusion system described by the Itô equations (w.r.t. independent Wiener processes  $w_t, W_t$ ) with a small parameter  $\varepsilon$ :

$$dX_t^\varepsilon = f(X_t^\varepsilon, Y_t^\varepsilon) dt + g(X_t^\varepsilon, Y_t^\varepsilon) dw_t, \quad X_0^\varepsilon = x_0, \quad (1.1)$$

$$\varepsilon dY_t^\varepsilon = F(Y_t^\varepsilon) + \sqrt{\varepsilon} G(Y_t^\varepsilon) dW_t, \quad Y_0^\varepsilon = y_0. \quad (1.2)$$

Hereafter,  $X_t^\varepsilon$  and  $Y_t^\varepsilon$  are referred as the “slow” and “fast” components respectively. All the functions  $f, g, F, G$ , entering in (1.1) and (1.2), are unknown and only the slow component  $X^\varepsilon$  is observed. The goal is to recover from the observations  $X_t^\varepsilon, 0 \leq t \leq T$ , some characteristics of the process  $X^\varepsilon$  which can be used for a further statistical analysis of this process or forecasting.

Examples of such problems meet, for instance, in satellite imaging, where  $X_t^\varepsilon$  describes the observed signal and  $Y_t^\varepsilon$  is used to describe rotation and vibration of the satellite. One more reasonable example is connected to asset price processes in financial markets. A weekly (or monthly) observed asset price process  $X^\varepsilon$  can be interpreted as



© 2000 Kluwer Academic Publishers. Printed in the Netherlands.

the “slow” component. If we are interested in some “global” (macro) characteristics of this process, then the influence of other components of the market can be modeled via the “fast” process  $Y_t^\varepsilon$ . Some other applications of such approach to the control theory can be found in Kushner (1990) or Liptser, Runggaldier, Taksar (1996).

Equations of the form  $dX_t = f(X_t + Y_t) dt + dw_t$  are often used to model regression problems with errors in regressors. It is well known, see e.g. Carrol and Hall (1988), Fan and Truong (1993) that the presence of the “error” component  $Y_t$  in the regressor variable makes the problem of estimating the regression function  $f$  much more difficult. Even if the distribution of  $Y_t$  is known, the optimal rate of estimating the function  $f$  is only logarithmic in the observation time. We do not assume special additive structure for the arguments of the drift function  $f$  and no information about the distribution of the noisy component  $Y$  is available. Instead we only assume that  $Y^\varepsilon$  is a fast oscillating process. We shall see that this qualitative assumption allows for a reasonable quality of estimation of the “averaged” drift function  $\bar{f}$  which describes the “macro” characteristics of the process  $X^\varepsilon$ .

It is well known from Khasminskii (1966) (see also Freidlin and Wentzell (1984), Veretennikov (1991)) that, under some regularity conditions on the functions  $F$  and  $G$  from (1.2),  $Y^\varepsilon$  is a fast oscillating ergodic process while the slow process  $X^\varepsilon$  obeys, so called, Bogolubov’s averaging principle. This roughly means that the distribution of the slow component is close to the distribution of the diffusion process  $X_t$  defined by the Itô equation

$$dX_t = \bar{f}(X_t) dt + \bar{g}(X_t) d\bar{w}_t, \quad (1.3)$$

where  $\bar{w}$  is some Wiener process and the drift and diffusion coefficients  $\bar{f}, \bar{g}$  are defined by averaging the original coefficients with respect to the stationary density  $p$  of the fast process:

$$\bar{f}(x) = \int f(x, y) p(y) dy \quad \text{and} \quad \bar{g}(x) = \left( \int g^2(x, y) p(y) dy \right)^{1/2}.$$

In other words, the “macro behavior” of the process  $X^\varepsilon$  is determined only by the averaged functions  $\bar{f}$  and  $\bar{g}$ . This naturally leads to the problem of statistical estimation of these functions from observations  $X_t^\varepsilon$ ,  $0 \leq t \leq T$ , where  $T$  is the *observation time*.

In this paper, we focus on estimating the *dynamic function*  $\bar{f}(x)$ . We do not discuss here the problem of estimating the diffusion coefficient  $\bar{g}$  since in the case of continuous observations, the required information about the function  $g$  can be exactly recovered from the data, Section 3.5 below. We also restrict ourselves to the problem of

pointwise estimation, that is, given a point  $x$ , we estimate the value  $\bar{f}(x)$ . We refer to Lepski, Mammen and Spokoiny (1997) for a discussion of the relation between pointwise and global estimation. Note that the problem of the pointwise estimation of the drift function  $f$  is closely connected to the problem of forecasting the process  $X^\varepsilon$ . Indeed, if we observe the process  $(X_t^\varepsilon)$  until the time-point  $T$ , and if we are interested in a behavior of the process in the nearest future after  $T$ , then we have to estimate  $\bar{f}(x)$  for  $x = X_t^\varepsilon$ .

The estimation theory for diffusion type processes is well developed under the parametric modeling when underlying functions (drift and diffusion) are specified up to a value of a finite dimensional parameter (cf. Kutoyants, 1984b, 1994) or Bhattacharya and Götze (1995). In contrast, nonparametric estimation is not studied in details. The known results concern only with statistical inference for diffusion models with a small noise or for ergodic diffusion and a large observation time  $T$ . Kutoyants (1984a) evaluated the minimax rate of estimation of the drift coefficient using a kernel type estimator. Genon-Catalot, Laredo and Picard (1992) applied wavelets. Locally polynomial estimators are described in Fan and Gijbels (1996). Milstein and Nussbaum (1994) established the LeCam equivalence between the diffusion model and the “white noise model”. Some pertinent results for autoregressive models in discrete time can be found in Doukhan and Ghindes (1980), Collomb and Doukhan (1983), Doukhan and Tsybakov (1993), Delyon and Juditsky (1997), Neumann (1998). A series of papers discusses simultaneous estimation of the drift and diffusion functions, among them Hall and Carroll (1989), Härdle and Tsybakov (1997), Ruppert et al (1997), Fan and Yao (1988).

In this paper, we assume neither ergodic properties of the slow component nor the large observation time  $T$ . This makes the problem more complicated. Additional difficulties come from the fact that the coefficients of the slow process are contaminated by the unobserved fast one. To our knowledge, nonparametric statistical inference for diffusion models (1.1), (1.2) with averaging has not yet been considered.

We propose a locally linear estimator of  $\bar{f}(x)$  with a data-driven bandwidth choice and show that this method provides a nearly optimal rate of estimation up to a log log factor.

The paper is organized as follows. The next section contains the description of the locally linear estimator. Its properties are discussed in Section 3. The data-driven bandwidth choice is presented in Section 4. All proofs are collected in Sections 5.

## 2. A locally linear estimator

For fixed  $x$ , to estimate the value  $\bar{f}(x)$  we apply the locally linear smoother (cf. Katkovnik (1985), Tsybakov (1986), Fan and Gijbels (1996)).

We begin with some heuristic explanations of the method. Imagine for a moment that the observed process  $X_t$ ,  $0 \leq t \leq T$  satisfies the Itô equation with respect to Wiener process  $w_t$ :

$$dX_t = f(X_t) dt + g(X_t) dw_t \quad (2.1)$$

with the linear function  $f$ :  $f(u) = \theta_0 + \theta_1 \frac{u-x}{h}$ , depending on two parameters  $\theta_0, \theta_1$ , where  $x$  and  $h > 0$  are fixed. These parameters can be estimated by the maximum likelihood method:

$$(\tilde{\theta}_0, \tilde{\theta}_1) = \operatorname{argmax}_{\theta_0, \theta_1} \left\{ \int_0^T \left( \theta_0 + \theta_1 \frac{X_t - x}{h} \right) dX_t - \frac{1}{2} \int_0^T \left( \theta_0 + \theta_1 \frac{X_t - x}{h} \right)^2 dt \right\},$$

that is, with  $\mu_k = \int_0^T \left( \frac{X_t - x}{h} \right)^k dt$ ,  $k = 0, 1, 2$ , we get

$$\tilde{\theta}_0 = \frac{\mu_2 \int_0^T dX_t - \mu_1 \int_0^T \frac{X_t - x}{h} dX_t}{\mu_0 \mu_2 - \mu_1^2},$$

$$\tilde{\theta}_1 = \frac{-\mu_1 \int_0^T dX_t + \mu_0 \int_0^T \frac{X_t - x}{h} dX_t}{\mu_0 \mu_2 - \mu_1^2}.$$

Since clearly  $f(x) = \theta_0$ , the value  $\tilde{\theta}_0$  can be taken as the estimate of  $f(x)$ .

The locally linear smoother is defined in a similar way. The only difference is that the function  $f$  is not assumed to be linear but it is approximated by a linear function  $\theta_0 + \theta_1 \frac{u-x}{h}$  in a small neighborhood  $[x-h, x+h]$  of the point  $x$ . Then the coefficients  $\theta_0, \theta_1$  of this function can be estimated from the observations of  $X_t$  falling into the interval  $[x-h, x+h]$ . For the formal description, let us introduce the *kernel* function  $K(u)$  which is assumed to be smooth, non-negative, bounded by 1, and vanishing outside of  $[-1, 1]$ . Then the locally linear estimate with the kernel  $K$  and a *bandwidth*  $h$  is defined as:

$$\tilde{f}_h(x) = \frac{\mu_{2,h} \int_0^T K\left(\frac{X_t - x}{h}\right) dX_t - \mu_{1,h} \int_0^T \frac{X_t - x}{h} K\left(\frac{X_t - x}{h}\right) dX_t}{\mu_{0,h} \mu_{2,h} - \mu_{1,h}^2}, \quad (2.2)$$

where

$$\mu_{k,h} = \int_0^T \left( \frac{X_t - x}{h} \right)^k K \left( \frac{X_t - x}{h} \right) dt, \quad k = 0, 1, 2.$$

Now we come back to the more complicated two-scaled model (1.1), (1.2). Here, due to the averaging principle, the observed process  $X_t^\varepsilon$  is closed in the distribution sense to the “limit” process  $X_t$  described by the equation (1.3). Therefore, to define our estimate  $\tilde{f}_h(x)$  of  $\tilde{f}(x)$ , we simply replace in the expression (2.2) the “limit” process  $X_t$  by our observations  $X_t^\varepsilon$ :

$$\tilde{f}_h(x) = \frac{\mu_{2,h} \int_0^T K \left( \frac{X_t^\varepsilon - x}{h} \right) dX_t^\varepsilon - \mu_{1,h} \int_0^T \frac{X_t^\varepsilon - x}{h} K \left( \frac{X_t^\varepsilon - x}{h} \right) dX_t^\varepsilon}{\mu_{0,h} \mu_{2,h} - \mu_{1,h}^2}, \quad (2.3)$$

where now

$$\mu_{k,h} = \int_0^T \left( \frac{X_t^\varepsilon - x}{h} \right)^k K \left( \frac{X_t^\varepsilon - x}{h} \right) dt, \quad k = 0, 1, 2. \quad (2.4)$$

The quality of estimate (2.3) essentially depends on the bandwidth  $h$ . Some useful properties of  $\tilde{f}_h(x)$  for the fixed  $h$  are described in Section 3. We discuss the adaptive choice of the bandwidth  $h$  in Section 4.

### 3. Accuracy of the locally linear estimate

In this section we study some properties of the locally linear estimate  $\tilde{f}_h(x)$  from (2.3). We first formulate the required conditions on the coefficients of the two-scaled system (1.1), (1.2). Then we present the result and discuss some its corollaries.

#### 3.1. CONDITIONS

In the sequel we suppose that the functions  $f, g$  and  $F, G$  from (1.1) and (1.2) obey the following conditions:

- (A<sub>s</sub>) Functions  $f(x, y)$  and  $g(x, y)$  are Lipschitz continuous in  $x, y$  and  $f(x, y)$  is three times continuously differentiable in  $x$ . For some positive constants  $g_{\min} \leq g_{\max}$

$$g_{\min} \leq |g(x, y)| \leq g_{\max}.$$

( $A_f$ ) 1. Functions  $F(y)$  and  $G(y)$  are Lipschitz continuous in  $y$  and continuously differentiable ( $F$  once,  $G$  twice) and their derivatives are continuous and bounded.

2. There exist constants  $\kappa > 0$  and  $C > 1$  such that for  $|y| > C$

$$yF(y) \leq -\kappa|y|^2,$$

3. Function  $G$  is bounded and strongly positive, i.e. for any  $y$

$$0 < G_{\min} \leq |G(y)| \leq G_{\max}.$$

Condition ( $A_f$ ) guarantees the required ergodicity of the fast process  $Y_t^\varepsilon$  and, moreover, this condition can be viewed as the mathematical formulation of the ergodic property of the fast process, see Veretennikov (1991) for more detailed analysis. Under ( $A_f$ ) the invariant density of the fast process can be explicitly described (Khasminskii, 1966) and it does not depend on  $\varepsilon$ :

$$p(y) = \text{Const.} \frac{\exp \left\{ 2 \int_0^y \frac{F(u)}{G^2(u)} du \right\}}{G^2(y)}. \quad (3.1)$$

It is worth to mention that neither the constants  $C, \kappa, G_{\min}, G_{\max}$ , nor the invariant density  $p$  are not assumed to be known and they do not enter into the description of the procedure and into the formulation of the main results.

### 3.2. ACCURACY OF THE LOCALLY LINEAR ESTIMATE

To state the result, we introduce some additional notations. With  $\mu_{k,h}$  defined in (2.4), set

$$D_h = \mu_{0,h}\mu_{2,h} - \mu_{1,h}^2, \quad (3.2)$$

and

$$\begin{aligned} \sigma_h^2(x) &= \frac{1}{D_h^2} \int_0^T \left( \mu_{2,h} - \mu_{1,h} \frac{X_t^\varepsilon - x}{h} \right)^2 K^2 \left( \frac{X_t^\varepsilon - x}{h} \right) g^2(X_t^\varepsilon, Y_t^\varepsilon) dt \\ &= v_{2,h}^2 V_{0,h} - 2v_{1,h}v_{2,h}V_{1,h} + v_{1,h}^2 V_{2,h} \end{aligned} \quad (3.3)$$

where

$$\begin{aligned} v_{k,h} &= \frac{\mu_{k,h}}{D_h} = \frac{\mu_{k,h}}{\mu_{0,h}\mu_{2,h} - \mu_{1,h}^2}, \quad k = 1, 2, \\ V_{k,h} &= \int_0^T \left( \frac{X_t^\varepsilon - x}{h} \right)^k K^2 \left( \frac{X_t^\varepsilon - x}{h} \right) g^2(X_t^\varepsilon, Y_t^\varepsilon) dt. \end{aligned}$$

Although the expressions for  $V_{k,h}$ ,  $k = 0, 1, 2$ , use the unknown diffusion coefficient  $g^2(X_t^\varepsilon, Y_t^\varepsilon)$  and moreover, one of its arguments  $Y_t^\varepsilon$  is not observed, these values can be computed on the base of our observations  $(X_t^\varepsilon, 0 \leq t \leq T)$  only, see Section 3.5.

The value  $\sigma_h^2(x)$  is called the *conditional variance* of the estimate  $\tilde{f}_h(x)$ . We use this terminology by analogy with the regression case, where  $X_t^\varepsilon$  is a deterministic design process and  $\sigma_h^2(x)$  is really the variance of the least squares estimate  $\tilde{f}_h(x)$ . Note that for the regression setup, some design regularity is required to ensure that  $\sigma_h^2(x)$  is not too large.

In our case,  $X_t^\varepsilon$  is the observed process which at the same time can be viewed as the design process. We therefore impose some conditions on the trajectories of the process  $X_t^\varepsilon$  which are similar to that of used to describe the design regularity in the regression setting. Our results are also similar to that of usually obtained in the regression estimation. In particular, we show that under the conditions imposed, the conditional variance  $\sigma_h^2(x)$  helps to control the stochastic component of the estimate  $\tilde{f}_h(x)$ .

For some  $\rho \geq 0$ ,  $r > 0$ ,  $b > 0$  and  $B \geq 1$  we introduce the set

$$\mathcal{A}_h = \left\{ \begin{array}{l} \frac{b}{Th} \leq v_{2,h} \leq \frac{bB}{Th}, \quad \frac{b}{Th} \leq \sigma_h^2(x) \leq \frac{bB}{Th}, \\ \mu_{0,h} \leq r\mu_{2,h}, \quad V_{0,h} \leq rV_{2,h}, \\ \mu_{1,h}^2 \leq \rho\mu_{0,h}\mu_{2,h}, \quad V_{1,h}^2 \leq \rho V_{0,h}V_{2,h} \end{array} \right\}.$$

Since  $X_t^\varepsilon$  is the random process, the set  $\mathcal{A}_h$  is random as well. In the sequel we study the properties of  $\tilde{f}_h(x)$  restricted to the set  $\mathcal{A}_h$ , see Section 3.3 for further discussion.

The quality of the approximation of  $f(u, y)$  by a linear in  $u$  function in the neighborhood  $u \in [x-h, x+h]$  is characterized by the following quantity

$$\Delta_h(x) = \sup_{|u-x| \leq h, y \in \mathbb{R}} |f(u, y) - f(x, y) - (u-x)f_x(x, y)|. \quad (3.4)$$

In the next theorem we describe some useful properties of the estimate (2.3).

**THEOREM 3.1.** *Let  $(A_s)$  and  $(A_f)$  be fulfilled, and let the values  $\varepsilon$  and  $\varepsilon T$  be sufficiently small and  $Th \geq 1$ . Then for every  $\lambda \geq \sqrt{2}$*

$$\mathbf{P} \left( \left| \tilde{f}_h(x) - \bar{f}(x) \right| > c\Delta_h(x) + \lambda\sigma_h(x), \mathcal{A}_h \right) \leq \alpha(\lambda) \quad (3.5)$$

with

$$\alpha(\lambda) = 4e \log(4B^3) \left( 1 + 4r \sqrt{\frac{1+r}{1-\rho}} \lambda^2 \right) \lambda e^{-\frac{\lambda^2}{2}} \quad (3.6)$$

and  $c = (1 - \rho)^{-1/2}$ .

Informally the result of the theorem means that for sufficiently large  $\lambda$ , the losses  $|\tilde{f}_h(x) - \bar{f}(x)|$  of the estimate  $\tilde{f}_h(x)$ , being restricted to  $\mathcal{A}_h$ , are bounded by the sum of two terms:  $c\Delta_h(x)$  and  $\lambda\sigma_h(x)$ . The first one mimics the accuracy of approximating the function  $f(u)$  by a linear in  $u$  function in the small vicinity  $[x-h, x+h]$  of  $x$ . The second term is in proportion to the “stochastic standard deviation”  $\sigma_h(x)$ .

### 3.3. SOME REMARKS RELATED TO THE RANDOM SET $\mathcal{A}_h$

The result of Theorem 3.1 describes the accuracy of the estimate  $\tilde{f}_h(x)$  on the random set  $\mathcal{A}_h$  only. Here we briefly discuss some related questions.

#### 3.3.1. Reason for restricting to $\mathcal{A}_h$

It was mentioned previously that restricting to  $\mathcal{A}_h$  allows to eliminate irregular cases when, for instance, the trajectory  $X_{[0,T]}^\varepsilon$  does not pass through the interval  $[x-h, x+h]$  and  $\mu_{0,h} = \mu_{1,h} = \mu_{2,h} = D_h = 0$ . Note that for typical applications to forecasting, one has to estimate  $\bar{f}(x)$  with  $x = X_t^\varepsilon$ , and the path  $X_{[0,T]}^\varepsilon$  obviously passes through  $x$ .

#### 3.3.2. Verifying the condition $X_{[0,T]}^\varepsilon \in \mathcal{A}_h$

Clearly the event  $\mathcal{A}_h$  is completely determined by the known values  $\mu_{k,h}$  and  $V_{k,h}$ ,  $k = 0, 1, 2$ . It is therefore always possible to check whether the observed trajectory  $X_{[0,T]}^\varepsilon$  belongs to  $\mathcal{A}_h$  or not. If the trajectory  $X_{[0,T]}^\varepsilon$  does not belong to  $\mathcal{A}_h$ , we are not able to guarantee a reasonable quality for the estimate  $\tilde{f}_h(x)$ .

#### 3.3.3. The conditions entering into the definition of $\mathcal{A}_h$

The conditions  $0 \leq K(u) \leq 1$  and  $K(u) = 0$  for  $|u| \geq 1$  imply  $\mu_{2,h} \leq \mu_{0,h}$  and  $V_{2,h} \leq V_{0,h}$ . Further, by the Cauchy-Schwarz inequality, it holds  $\mu_{1,h}^2 \leq \mu_{0,h}\mu_{2,h}$  and  $V_{1,h}^2 \leq V_{0,h}V_{2,h}$ . The conditions  $\mu_{0,h} \leq r\mu_{2,h}$ ,  $V_{0,h} \leq rV_{2,h}$ ,  $\mu_{1,h}^2 \leq \rho\mu_{0,h}\mu_{2,h}$  and  $V_{1,h}^2 \leq \rho V_{0,h}V_{2,h}$  with  $\rho < 1$  and  $r \geq 1$  ensure that the local linear estimate is well defined. Note that these conditions are not completely independent. In particular, if  $g(x)$  is a constant function and if  $K(u) = 1(|u| \leq 1)$ , then  $\mu_{k,h} = V_{k,h}$  for  $k = 0, 1, 2$  and  $\sigma_h^2(x) = v_{2,h} = \mu_{2,h}/(\mu_{0,h}\mu_{2,h} - \mu_{1,h}^2)$ .



### 3.3.4. The choice of the constants $\rho$ , $b$ , $B$ , $r$

The choice of constants  $\rho$ ,  $b$ ,  $B$ ,  $r$ , entering in the definition of the set  $\mathcal{A}_h$ , is optional and they even may depend on  $T$ .

For a regular design in the regression setup, it holds  $\mu_{1,h} = V_{1,h} = 0$ . If, in addition,  $g(u)$  is constant in the interval  $[x - h, x + h]$ , then  $\mu_{0,h} = r(K)\mu_{2,h}$  and  $V_{0,h} = r(K)V_{2,h}$  with

$$r(K) = \int K(u) du \left( \int u^2 K(u) du \right)^{-1}.$$

Therefore, I reasonable choice would be  $\rho = 1/2$  and  $r = 2r(K)$ .

Concerning the choice of the parameters  $b, B$ , note that the upper bound (3.5) from Theorem 3.1 does not depends on  $b$  and it depends on  $B$  (which determines the range of different values for the conditional variance  $\sigma_h^2(x)$ ) only via the log-factor  $\log(4B^3)$ . Simple heuristic consideration prompt a possible choice  $b = h_{\min}$  and  $B = T$ .

### 3.3.5. Unconditional result under ergodicity

If the coefficients  $f$  and  $g$  obey some additional conditions which ensure ergodicity of the process  $X_t^\varepsilon$ , see e.g. Veretennikov (1991), then, at least with growing  $T$  the normalized integrals  $(Th)^{-1}\mu_{k,h}$  and  $(Th)^{-1}V_{k,h}$  ( $k = 0, 1, 2$ ) converge to some fixed values which depend only on the stationary distribution of the process  $X_t^\varepsilon$ . Moreover, one can usually select fixed constants  $b, B$  and  $\rho, r$  in such a way that  $1 - \mathbf{P}(\mathcal{A}_h)$  converges to zero exponentially fast as  $T \rightarrow \infty$ . Since obviously

$$\begin{aligned} \mathbf{P} \left( \left| \tilde{f}_h(x) - \bar{f}(x) \right| > c\Delta_h(x) + \lambda\sigma_h(x) \right) \\ \leq \mathbf{P} \left( \left| \tilde{f}_h(x) - \bar{f}(x) \right| > c\Delta_h(x) + \lambda\sigma_h(x), \mathcal{A}_h \right) + \mathbf{P}(\mathcal{A}_h^c) \end{aligned}$$

we obtain in this situation an unconditional asymptotic bound for the risk of the estimate  $\tilde{f}_h(x)$ .

## 3.4. QUALITY OF ESTIMATION UNDER SMOOTHNESS ASSUMPTIONS

Due to the assumptions  $(A_s)$  from Section 3, the function  $f$  is twice continuously differentiable with respect to the first argument. Assume also that for every  $u$  from a small vicinity of  $x$  and any fixed  $y$

$$\left| \frac{\partial^2 f(u, y)}{\partial u^2} \right| \leq L. \quad (3.7)$$

Then the value  $\Delta_h(x)$  defined in (3.4), is bounded above by  $Lh^2/2$ . On the other hand, on the set  $\mathcal{A}_h$  the stochastic variance  $\sigma_h^2(x)$  is of

order  $(Th)^{-1}$ . Therefore, following to the standard approach in non-parametric estimation, the bandwidth  $h$  can be chosen by balancing the accuracy of approximation and the stochastic error:

$$Lh^2 \asymp \frac{1}{\sqrt{T}h}.$$

This leads to the choice  $h \asymp (TL^2)^{-1/5}$  and hence to the rate of the estimation  $L^{1/5}T^{-2/5}$  which is optimal in the minimax sense under the smoothness assumptions (3.7), see e.g. Ibragimov and Khasmiskii (1981). Unfortunately this approach hardly applies in practice, since the constant  $L$  in (3.7) is typically unknown. An adaptive (data-driven) choice of the bandwidth is discussed in the next section.

### 3.5. COMPUTATION OF $\sigma_h^2(x)$

Recall that with fixed  $h$ , the value  $\sigma_h^2(x)$  is defined by the formula

$$\begin{aligned} \sigma_h^2(x) &= \frac{1}{D_h^2} \int_0^T K^2 \left( \frac{X_t^\varepsilon - x}{h} \right) \left( \mu_{2,h} - \mu_{1,h} \frac{X_t^\varepsilon - x}{h} \right)^2 g^2(X_t^\varepsilon, Y_t^\varepsilon) dt \\ &= v_{2,h}^2 V_{0,h} - 2v_{1,h}v_{2,h}V_{1,h} + v_{1,h}^2 V_{2,h} \end{aligned}$$

with

$$\begin{aligned} \mu_{k,h} &= \int_0^T \left( \frac{X_t^\varepsilon - x}{h} \right)^k K \left( \frac{X_t^\varepsilon - x}{h} \right) dt, \\ D_h &= \mu_{0,h}\mu_{2,h} - \mu_{1,h}^2, \\ v_{k,h} &= \frac{\mu_{k,h}}{D_h} = \frac{\mu_{k,h}}{\mu_{0,h}\mu_{2,h} - \mu_{1,h}^2}, \\ V_{k,h} &= \int_0^T \left( \frac{X_t^\varepsilon - x}{h} \right)^k K^2 \left( \frac{X_t^\varepsilon - x}{h} \right) g^2(X_t^\varepsilon, Y_t^\varepsilon) dt, \end{aligned}$$

for  $k = 0, 1, 2$ . The formula for  $\sigma_h^2(x)$  includes the unknown diffusion coefficient  $g^2(X_t^\varepsilon, Y_t^\varepsilon)$  and the unobserved process  $Y_t^\varepsilon$  as one of its arguments. We now show that despite of this fact, the value  $\sigma_h^2(x)$  can be computed via the trajectory  $X_{[0,T]}^\varepsilon$  only.

Let us introduce two random processes

$$Z_t' = \int_0^t K \left( \frac{X_s^\varepsilon - x}{h} \right) dX_s^\varepsilon \quad \text{and} \quad Z_t'' = \int_0^t K \left( \frac{X_s^\varepsilon - x}{h} \right) \frac{X_s^\varepsilon - x}{h} dX_s^\varepsilon$$

which are completely determined on the time interval  $[0, T]$  by  $X_{[0,T]}^\varepsilon$ . Applying the Itô formula we get

$$(Z_T')^2 = 2 \int_0^T Z_t' dZ_t' + V_{0,h}$$

$$\begin{aligned}(Z_T'')^2 &= 2 \int_0^T Z_t'' dZ_t'' + V_{2,h} \\ Z_T' Z_T'' &= \int_0^T Z_t' dZ_t'' + \int_0^T Z_t'' dZ_t' + V_{1,h}.\end{aligned}$$

Hence  $V_{0,h} = (Z_T')^2 - 2 \int_0^T Z_t' dZ_t'$ , so that  $V_{0,h}$  is completely determined by  $X_{[0,T]}^\varepsilon$ . Similar arguments apply for  $V_{1,h}$  and  $V_{2,h}$  and hence for  $\sigma_h^2(x)$  as required.

#### 4. Data-driven bandwidth selection

In this section we consider the problem of bandwidth selection for the locally linear estimator described in Section 2. It is assumed here that the method of estimation, that is, the locally linear smoother with the kernel  $K$ , is fixed and only the bandwidth  $h$  has to be chosen. Below we discuss one adaptive (data driven) approach which goes back to the idea of pointwise adaptive estimation, see Lepski (1990), Lepski and Spokoiny (1997) and Spokoiny (1998).

The idea of the method can be explained as follows. In the light of Theorem 3.1, we would be interested to select a bandwidth  $h$  which leads to a possibly small sum of the form  $c\Delta_h(x) + \lambda\sigma_h(x)$  among all considered bandwidth values  $h$ . This sum is comprised of two terms. The first one (“bias”) characterizes the accuracy of local approximation of the underlying drift function  $f$  by the linear functions and it typically increases with  $h$ . The second term is proportional to the conditional standard deviation  $\sigma_h(x)$  which typically decreases with  $h$ . (Indeed, an increase of  $h$  makes the estimation window  $[x-h, x+h]$  larger and hence more observations can be used for estimating the underlying function  $f$  at the point  $x$ . This results in a smaller variance of the estimate.) To simplify the exposition, we suppose that  $\sigma_h^2(x)$  strongly decreases in  $h \in \mathcal{H}$ . (If this assumption is not fulfilled for the original set  $\mathcal{H}$ , i.e. if there is  $h' < h \in \mathcal{H}$  with the property  $\sigma_h^2(x) \geq \sigma_{h'}^2(x)$ , then we simply exclude  $h$  from  $\mathcal{H}$ .)

Therefore, a “good” (or “ideal”) choice  $h_{\text{id}}$  corresponds to a possibly large bandwidth  $h$  (which makes the stochastic component of the estimate small) still providing that the “bias” component  $c\Delta_h(x)$  is not significantly larger than  $\sigma_h(x)$ . (We call  $h_{\text{id}}$  an “ideal” bandwidth since its definition relies on the unknown function  $\Delta_h(x)$ .) The latter property is clearly fulfilled for all smaller bandwidths  $h \leq h_{\text{id}}$ . Therefore, if  $h_{\text{id}}$  is “good” and  $h < h_{\text{id}}$ , then the two corresponding estimates  $\hat{f}_{h_{\text{id}}}(x)$  and  $\hat{f}_h(x)$  should not differ significantly.

The proposed procedure can be viewed as a family of tests whether the estimate  $\tilde{f}_h(x)$  for a bandwidth-candidate  $h$  does not differ significantly from estimates  $\tilde{f}_\eta(x)$  with smaller bandwidths  $\eta < h$ . The latter is done on the base of Theorem 3.1 which allows to bound with a large probability the difference  $|\tilde{f}_h(x) - \tilde{f}_\eta(x)|$  by  $\lambda\sigma_h(x) + \lambda\sigma_\eta(x) + c\Delta_h(x) + c\Delta_\eta(x)$  provided that  $\lambda$  is sufficiently large. The terms  $c\Delta_h(x)$  and  $c\Delta_\eta(x)$  in this sum are unknown but, if  $h$  is “good” that is, if  $\Delta_h(x) \ll \sigma_h(x)$ , then their contribution is negligible. In opposite, a significant deviation of  $|\tilde{f}_h(x) - \tilde{f}_\eta(x)|$  over the level  $\lambda\sigma_h(x) + \lambda\sigma_\eta(x)$  can be explained only by a large bias component indicating that  $h$  is not a “good” bandwidth. The procedure searches for the largest bandwidth  $h$  such that the hypothesis  $\tilde{f}_h(x) = \tilde{f}_\eta(x)$  is not rejected for all  $\eta < h$ .

Now we present a formal description. Suppose a family  $\mathcal{H}$  of bandwidth-candidates  $h$  is fixed. For technical reasons, we assume that this set is finite and denote by  $H$  the number of its elements. Usually  $\mathcal{H}$  is taken as a geometric grid:

$$\mathcal{H} = \{h = h_{\min}a^k, k = 0, 1, 2, \dots : h \leq h_{\max}\},$$

where  $h_{\min} \leq h_{\max}$  and  $a > 1$  are some prescribed constants. As in Section 3 we restrict ourselves only to those  $h$  from  $\mathcal{H}$  for which the observed path  $X_{[0,T]}^\varepsilon$  belongs to  $\mathcal{A}_h$ .

With every bandwidth value  $h$  we associate the estimate  $\tilde{f}_h(x)$  of  $\bar{f}(x)$  and the corresponding conditional standard deviations  $\sigma_h(x)$  which can be precisely calculated as described in Section 3.5.

Now, with two constants  $\lambda_1$  and  $\lambda_2$ , define the adaptive choice of bandwidth by the following iterative procedure:

**Initialization** Select the smallest bandwidth in  $\mathcal{H}$ .

**Iteration** Select the next larger bandwidth  $h$  in  $\mathcal{H}$  and calculate the corresponding estimate  $\tilde{f}_h(x)$  and the conditional standard deviation  $\sigma_h(x)$ .

**Testing** Reject  $h$ , if there exists one  $\eta \in \mathcal{H}$  with  $\eta < h$  such that

$$|\tilde{f}_h(x) - \tilde{f}_\eta(x)| > \lambda_1 \sigma_\eta(x) + \lambda_2 \sigma_h(x). \quad (4.1)$$

**Loop** If  $h$  is not rejected, then continue with *iteration step* by choosing a larger bandwidth  $h$  in  $\mathcal{H}$ . Otherwise, set  $\hat{h} =$  ”the latest non rejected  $h$ ”.

The proposed rule can be packed in the following form:

$$\hat{h} = \max \left\{ h \in \mathcal{H} : |\tilde{f}_{h'}(x) - \tilde{f}_\eta(x)| \leq \lambda_1 \sigma_\eta(x) + \lambda_2 \sigma_{h'}(x) \right\} \quad (4.2)$$

$$\left. \forall h', \eta \in \mathcal{H}, \eta < h' \leq h \right\}.$$

The choice of the parameters  $\lambda_1, \lambda_2$  and the set  $\mathcal{H}$  is discussed in Section 4.1.

Finally, to define our adaptive estimate, we plug the data-driven bandwidth  $\hat{h}$  in the estimate  $\tilde{f}_h(x)$ , that is,  $\hat{f}(x) \equiv \tilde{f}_{\hat{h}}(x)$ .

In the next theorem we describe some properties of the adaptive estimate  $\hat{f}(x)$  restricted to the set

$$\mathcal{A}^* = \bigcap_{h \in \mathcal{H}} \mathcal{A}_h.$$

**THEOREM 4.1.** *The estimate  $\hat{f}(x) \equiv \tilde{f}_{\hat{h}}(x)$  with  $\hat{h}$  from (4.2) and  $\lambda_2 \geq \lambda_1$  fulfills the following property:*

$$\mathbf{P} \left( \left| \hat{f}(x) - \bar{f}(x) \right| > (2\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x), \mathcal{A}^* \right) \leq \sum_{\eta \in \mathcal{H} : \eta \leq h_{\text{id}}} \alpha(\lambda_\eta) \quad (4.3)$$

where  $\alpha(\lambda)$  is defined in (3.6) and

$$\lambda_\eta = \lambda_1 - c\Delta_\eta(x)/\sigma_\eta(x). \quad (4.4)$$

#### 4.1. THE CHOICE OF PARAMETERS $\lambda_1, \lambda_2, h_{\min}, h_{\max}$ AND $a$

Different proposals for the choice of the grid  $\mathcal{H}$  is discussed in Lepski, Mammen and Spokoiny (1997) and in Lepski and Spokoiny (1997). One possible choice for the grid  $\mathcal{H}$  reads as follows:  $h_{\min} = 1/T$ ,  $h_{\max} = 1$ ,  $a = \sqrt{2}$ , although these values can be changed without essential influence on the quality of the procedure.

The choice of parameters  $\lambda_1, \lambda_2$ , entering in (4.2), plays more important role. We start with the following general remark: the upper bound for the risk from Theorem 4.1 is rather rough and should be used with care for the parameter selection. However, it delivers some useful qualitative information about this choice which can be used for a theoretical study. The bound in (4.3) shows that the probability for  $|\hat{f}(x) - \bar{f}(x)|$  of being large is small, provided that the value  $\sum_{\eta \in \mathcal{H} : \eta \leq h_{\text{id}}} \alpha(\lambda_\eta)$  is sufficiently small. Here we discuss shortly the specific case when the values  $\Delta_\eta(x)$  vanish. The general case can be relatively easily reduced to that one. Indeed, a “good” bandwidth  $h_{\text{id}}$  can be defined by trade-off arguments between the “bias”  $c\Delta_{h_{\text{id}}}(x)$  and the conditional standard deviation  $\sigma_{h_{\text{id}}}(x)$ , that is,  $h_{\text{id}}$  is the maximal  $h$  from  $\mathcal{H}$  with  $c\Delta_h(x) \leq D\sigma_h(x)$  for some fixed value  $D$ . Taking  $D$  small enough provides that  $c\Delta_\eta(x) \ll \sigma_\eta(x)$  for all  $\eta \leq h_{\text{id}}$ .

If  $\Delta_\eta(x)$  vanishes for all such  $\eta$ , then  $\lambda_\eta = \lambda_1$  and

$$\sum_{\eta \in \mathcal{H} : \eta \leq h_{\text{id}}} \alpha(\lambda_\eta) \leq H\alpha(\lambda_1).$$

Therefore,  $\lambda_1$  should be selected in a way to provide that  $H\alpha(\lambda_1)$  is sufficiently small. This leads to the choice

$$\lambda_1 \approx \sqrt{2 \log(H) + \lambda^2}$$

with some fixed constant  $\lambda$  so that

$$He^{-\lambda_1^2/2} \approx e^{-\lambda^2/2}.$$

If  $\mathcal{H}$  is taken in the form of the geometric grid, then we get  $H \approx \log_a(h_{\text{max}}/h_{\text{min}})$ . Therefore, taking  $h_{\text{max}} \approx 1$  and  $h_{\text{min}} \approx 1/T$ , we arrive at

$$\lambda_1 \approx \sqrt{2 \log \log T + \lambda^2}.$$

There is much more degree of freedom in the choice of  $\lambda_2$ . The constraint  $\lambda_2 \geq \lambda_1$  from Theorem 4.1 is of technical matter and it is used only in theoretical investigations. It can be skipped in practical applications. Simulation results show a reasonable (and very similar) performance of the presented procedure with  $\lambda_1 \approx 2$  and  $\lambda_2 = 1$ , or  $\lambda_1 = \lambda_2 = 1.5$  in the most cases. We refer to the forthcoming paper by Mercurio and Spokoiny (2000) for a more detailed discussion of practical issues and for a proposal of a data-driven choice of the parameters  $\lambda_1$  and  $\lambda_2$  in the context of applications to finance time series.

#### 4.2. ACCURACY OF ADAPTIVE ESTIMATION

We now compare the accuracy of the adaptive procedure (4.2) with the “optimal” one designed for the case of known smoothness properties of the underlying function  $f$  (see Section 3.4).

Assume  $|f''(u)| \leq L$ , see (3.7). Then  $\Delta_h(x) \leq Lh^2/2$  and the conditions  $\sigma_h^2(x) \asymp (hT)^{-1}$  and the balance relation  $c\Delta_h(x) \leq D\sigma_h(x)$  yield for  $h_{\text{id}}$ :

$$h_{\text{id}} \asymp (TL^2)^{-1/5}$$

so that  $\sigma_{h_{\text{id}}}(x) \asymp L^{1/5}T^{-2/5}$ . Hence, the above-mentioned choice  $\lambda_1 \approx \sqrt{2 \log \log T}$  and  $\lambda_2 = \lambda_1$  leads due to Theorem 4.1 to the following accuracy of the adaptive estimation:

$$(2\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x) \asymp L^{1/5} \left( \frac{\log \log T}{T} \right)^{2/5}.$$

At the same time, the “ideal” choice of the bandwidth leads to the rate  $L^{1/5}T^{-2/5}$ , see Section 3.4. Thus, the accuracy of adaptive estimation is worse than the “ideal” one within a  $\log \log T$ -factor only.

The origin of the  $\log \log T$ -factor in the rate of adaptive estimation can be easily explained. The total number  $H$  of considered estimates is logarithmic in the observation time  $T$  and the adaptive choice of the bandwidth leads to a worse accuracy by factor  $\log(H)$  at some power.

The notion of “payment for adaptation” is now well understood in nonparametric estimation: if we have too many estimates to select between, we have to “pay” for the adaptive choice some additional factor in the risk of estimation. In particular, it is shown in Lepski (1990) and Brown and Low (1996) (see also Lepski and Spokoiny (1997)) that for the problem of pointwise adaptive estimation, the optimal adaptive rate has to be worse than the optimal one by a log-factor.

In our results a  $\log \log$ -factor appears. This fact is not in contradiction to earlier issues, since the above-mentioned results correspond to the case of the power loss function  $\ell(x) = |x|^p$ ,  $p > 0$ , while we consider the bounded loss function. It can be also shown that the rate achieved by our estimate is optimal for pointwise adaptive estimation with a bounded loss function (see Spokoiny (1997) for similar results in the adaptive testing problem).

## 5. Proofs

In this section we prove Theorems 3.1 and 4.1. For a generic positive constant the notation ‘ $\ell$ ’ will be used hereafter.

### 5.1. DECOMPOSITION OF $\tilde{f}_h(x)$

We use two obvious identities characterizing the local linear smoother: for  $v_{1,h} = \frac{\mu_{1,h}}{D_h}$  and  $v_{2,h} = \frac{\mu_{2,h}}{D_h}$

$$\int_0^T K\left(\frac{X_s^\varepsilon - x}{h}\right) \left(v_{2,h} - v_{1,h} \frac{X_s^\varepsilon - x}{h}\right) ds = 1$$

$$\int_0^T K\left(\frac{X_s^\varepsilon - x}{h}\right) \left(v_{2,h} \frac{X_s^\varepsilon - x}{h} - v_{1,h} \frac{(X_s^\varepsilon - x)^2}{h^2}\right) ds = 0$$

and hence, with  $U_{s,h}^\varepsilon = \frac{X_s^\varepsilon - x}{h}$ ,

$$\int_0^T K(U_{s,h}^\varepsilon) (v_{2,h} - v_{1,h} U_{s,h}^\varepsilon) \bar{f}(x) ds = \bar{f}(x) \quad (5.1)$$

$$\int_0^T K(U_{s,h}^\varepsilon) (v_{2,h} U_{s,h}^\varepsilon - v_{1,h} (U_{s,h}^\varepsilon)^2) \bar{f}_x(x) ds = 0. \quad (5.2)$$

Due to (2.3) and (1.1), the estimate  $\tilde{f}_h(x)$  can be represented as follows:

$$\begin{aligned}
\tilde{f}_h(x) &= v_{2,h} \int_0^T K(U_{s,h}^\varepsilon) dX_s^\varepsilon \\
&\quad - v_{1,h} \int_0^T K(U_{s,h}^\varepsilon) U_{s,h}^\varepsilon dX_s^\varepsilon \\
&= \int_0^T K(U_{s,h}^\varepsilon) \left( v_{2,h} - v_{1,h} \frac{X_s^\varepsilon - x}{h} \right) f(X_s^\varepsilon, Y_s^\varepsilon) ds \\
&\quad + v_{2,h} \int_0^T K(U_{s,h}^\varepsilon) g(X_s^\varepsilon, Y_s^\varepsilon) dw_s \\
&\quad - v_{1,h} \int_0^T K(U_{s,h}^\varepsilon) U_{s,h}^\varepsilon g(X_s^\varepsilon, Y_s^\varepsilon) dw_s.
\end{aligned}$$

Now (5.1) and (5.2) imply the following decomposition

$$\tilde{f}_h(x) = \bar{f}(x) + \xi_h + r_h + \zeta_h^{(1)} + \zeta_h^{(2)} \quad (5.3)$$

where, with  $\delta(X_s^\varepsilon, Y_s^\varepsilon, x) = f(X_s^\varepsilon, Y_s^\varepsilon) - f(x, Y_s^\varepsilon) - \frac{X_s^\varepsilon - x}{h} f_x(x, Y_s^\varepsilon)$ ,

$$\begin{aligned}
r_h &= \int_0^T K(U_{s,h}^\varepsilon) (v_{2,h} - v_{1,h} U_{s,h}^\varepsilon) \delta(X_s^\varepsilon, Y_s^\varepsilon, x) ds, \\
\xi_h &= v_{2,h} \int_0^T K(U_{s,h}^\varepsilon) g(X_s^\varepsilon, Y_s^\varepsilon) dw_s \\
&\quad - v_{1,h} \int_0^T K(U_{s,h}^\varepsilon) U_{s,h}^\varepsilon g(X_s^\varepsilon, Y_s^\varepsilon) dw_s, \\
\zeta_h^{(1)} &= v_{2,h} \int_0^T K(U_{s,h}^\varepsilon) [f(x, Y_s^\varepsilon) - \bar{f}(x)] ds \\
&\quad - v_{1,h} \int_0^T K(U_{s,h}^\varepsilon) [f(x, Y_s^\varepsilon) - \bar{f}(x)] U_{s,h}^\varepsilon ds, \\
\zeta_h^{(2)} &= v_{2,h} \int_0^T K(U_{s,h}^\varepsilon) [f_x(x, Y_s^\varepsilon) - \bar{f}_x(x)] U_{s,h}^\varepsilon ds \\
&\quad - v_{1,h} \int_0^T K(U_{s,h}^\varepsilon) [f_x(x, Y_s^\varepsilon) - \bar{f}_x(x)] \frac{(X_s^\varepsilon - x)^2}{h^2} ds.
\end{aligned}$$

Below we evaluate separately each term in this decomposition.



5.2. AN UPPER BOUND FOR  $|r_h|$ 

Since  $K\left(\frac{u-x}{h}\right)$  vanishes for any  $u \notin [x-h, x+h]$  and  $|\delta(X_s^\varepsilon, Y_s^\varepsilon, x)| \leq \Delta_h(x)$  for  $|X_s^\varepsilon - x| \leq h$ , we get

$$\begin{aligned} |r_h| &\leq \int_0^T K(U_{s,h}^\varepsilon) (v_{2,h} - v_{1,h} U_{s,h}^\varepsilon) |\delta(X_s^\varepsilon, Y_s^\varepsilon, x)| \, ds \\ &\leq \Delta_h(x) \int_0^T K(U_{s,h}^\varepsilon) \left| v_{2,h} - v_{1,h} \frac{X_s^\varepsilon - x}{h} \right| \, ds. \end{aligned} \quad (5.4)$$

The properties  $|K(u)| \leq 1$  and  $K(u) = 0$ ,  $|u| \geq 1$  imply the inequality  $\mu_{2,h} \leq \mu_{0,h}$ . In addition we know that it holds on  $\mathcal{A}_h$

$$\mu_{1,h}^2 \leq \rho \mu_{0,h} \mu_{2,h}. \quad (5.5)$$

We now show that

$$|r_h| \leq (1 - \rho)^{-1/2} \Delta_h(x) \quad \text{on } \mathcal{A}_h. \quad (5.6)$$

The Cauchy-Schwarz inequality applied to (5.4) gives

$$r_h^2 \leq \Delta_h^2(x) \int_0^T K(U_{s,h}^\varepsilon) \, ds \int_0^T K(U_{s,h}^\varepsilon) (v_{2,h} - v_{1,h} U_{s,h}^\varepsilon)^2 \, ds.$$

Next,

$$\int_0^T K(U_{s,h}^\varepsilon) \, ds = \mu_{0,h},$$

and using  $v_{k,h} = \mu_{k,h}/D_h$ , with  $D_h = \mu_{2,h}\mu_{0,h} - \mu_{1,h}^2$ ,  $k = 0, 1, 2$ , we get

$$\begin{aligned} &\int_0^T K(U_{s,h}^\varepsilon) (v_{2,h} - v_{1,h} U_{s,h}^\varepsilon)^2 \, ds \\ &= \frac{1}{D_h^2} \int_0^T K(U_{s,h}^\varepsilon) (\mu_{2,h} - \mu_{1,h} U_{s,h}^\varepsilon)^2 \, ds \\ &= \frac{\mu_{2,h}^2}{D_h^2} \int_0^T K(U_{s,h}^\varepsilon) \, ds + \frac{\mu_{1,h}^2}{D_h^2} \int_0^T K(U_{s,h}^\varepsilon) (U_{s,h}^\varepsilon)^2 \, ds \\ &\quad - \frac{2\mu_{1,h}\mu_{2,h}}{D_h^2} \int_0^T K(U_{s,h}^\varepsilon) U_{s,h}^\varepsilon \, ds \\ &= \frac{\mu_{2,h}^2 \mu_{0,h} - \mu_{2,h} \mu_{1,h}^2}{D_h^2} \\ &= \mu_{2,h}/D_h. \end{aligned}$$

Hence, in view of (5.5),

$$r_h^2 \leq \Delta_h^2(x) \frac{\mu_{0,h} \mu_{2,h}}{D_h} = \Delta_h^2(x) \frac{\mu_{0,h} \mu_{2,h}}{\mu_{0,h} \mu_{2,h} - \mu_{1,h}^2} \leq \Delta_h^2(x) \frac{1}{1 - \rho}$$

as required.

### 5.3. AN UPPER BOUND FOR $\xi_h$

We study here some properties of the “stochastic term”

$$\begin{aligned} \xi_h &= v_{2,h} \int_0^T K \left( \frac{X_s^\varepsilon - x}{h} \right) g(X_s^\varepsilon, Y_s^\varepsilon) dw_s \\ &\quad - v_{1,h} \int_0^T K \left( \frac{X_s^\varepsilon - x}{h} \right) \frac{X_s^\varepsilon - x}{h} g(X_s^\varepsilon, Y_s^\varepsilon) dw_s. \end{aligned}$$

Namely, we intend to show that the probability of the event  $\{\xi_h > \lambda \sigma_h(x)\}$  with  $\sigma_h(x)$  from (3.3) is small provided that  $\lambda$  is large enough. Set for  $t \leq T$

$$\begin{aligned} M_{0,t} &= \int_0^t K \left( \frac{X_s^\varepsilon - x}{h} \right) g(X_s^\varepsilon, Y_t^\varepsilon) dw_s, \\ M_{1,t} &= \int_0^t K \left( \frac{X_s^\varepsilon - x}{h} \right) \frac{X_s^\varepsilon - x}{h} g(X_s^\varepsilon, Y_t^\varepsilon) dw_s. \end{aligned}$$

The Itô integrals  $M_{0,t}$  and  $M_{1,t}$  are continuous local martingales with the predictable quadratic variations (see e.g. Liptser and Shiryaev (1989))

$$\begin{aligned} \langle M_0 \rangle_t &= \int_0^t K^2 \left( \frac{X_s^\varepsilon - x}{h} \right) g^2(X_s^\varepsilon, Y_s^\varepsilon) ds, \\ \langle M_0, M_1 \rangle_t &= \int_0^t K^2 \left( \frac{X_s^\varepsilon - x}{h} \right) \frac{X_s^\varepsilon - x}{h} g^2(X_s^\varepsilon, Y_s^\varepsilon) ds, \\ \langle M_1 \rangle_t &= \int_0^t K^2 \left( \frac{X_s^\varepsilon - x}{h} \right) \left( \frac{X_s^\varepsilon - x}{h} \right)^2 g^2(X_s^\varepsilon, Y_s^\varepsilon) ds, \end{aligned}$$

so that  $\langle M_0 \rangle_T = V_{0,h}$ ,  $\langle M_0, M_1 \rangle_T = V_{1,h}$  and  $\langle M_1 \rangle_T = V_{2,h}$ . This yields

$$\begin{aligned} \xi_h(x) &= v_{2,h} M_{0,T} - v_{1,h} M_{1,T}, \\ \sigma_h^2(x) &= v_{2,h}^2 \langle M_0 \rangle_T - 2v_{1,h} v_{2,h} \langle M_0, M_1 \rangle_T + v_{1,h}^2 \langle M_1 \rangle_T. \end{aligned}$$

Denote

$$u_h = \frac{v_{1,h}}{v_{2,h}} = \frac{\mu_{1,h}}{\mu_{2,h}}.$$

Obviously

$$\begin{aligned} & \mathbf{P}(|\xi_h| > \lambda\sigma_h(x), \mathcal{A}_h) \\ &= \mathbf{P}\left(|M_{0,T} - u_h M_{1,T}| > \lambda\sqrt{V_T(u_h)}, \mathcal{A}_h\right). \end{aligned}$$

with  $V_T(u_h) = \langle M_0 \rangle_T - 2u_h \langle M_0, M_1 \rangle_T + u_h^2 \langle M_1 \rangle_T$ . To evaluate from above the right side of this equality, we apply the general result from Proposition 6.2, see Appendix. First we check the required conditions. The value  $|u_h|$ , being restricted to  $\mathcal{A}_h$ , can be bounded as:

$$|u_h| \leq \left| \frac{\sqrt{\rho\mu_{0,h}\mu_{2,h}}}{\mu_{2,h}} \right| \leq \sqrt{\rho r}.$$

Note now that

$$\begin{aligned} & \frac{\langle M_1 \rangle_T}{\langle M_0 \rangle_T - 2u_h \langle M_0, M_1 \rangle_T + u_h^2 \langle M_1 \rangle_T} \\ &= \frac{V_{2,h}}{V_{0,h} - 2u_h V_{1,h} + u_h^2 V_{2,h}} \\ &= \frac{V_{2,h}^2}{V_{0,h} V_{2,h} - V_{1,h}^2 + (V_{1,h} - u_h V_{2,h})^2}, \end{aligned}$$

and it holds on  $\mathcal{A}_h$  in view of  $V_{2,h} \leq V_{0,h}$

$$\frac{\langle M_1 \rangle_T}{\langle M_0 \rangle_T - 2u_h \langle M_0, M_1 \rangle_T + u_h^2 \langle M_1 \rangle_T} \leq \frac{V_{2,h}^2}{(1-\rho)V_{0,h}V_{2,h}} \leq \frac{1}{1-\rho}.$$

In addition, the definition of  $\mathcal{A}_h$  provides the following bounds for  $\sigma_h^2(x)$  on this set

$$\begin{aligned} \frac{\sigma_h^2(x)}{Th v_{2,h}^2} &= \frac{Th \sigma_h^2(x)}{(Th v_{2,h})^2} \leq \frac{bB}{b^2} = \frac{B}{b}, \\ \frac{\sigma_h^2(x)}{Th v_{2,h}^2} &= \frac{Th \sigma_h^2(x)}{(Th v_{2,h})^2} \geq \frac{b}{(bB)^2} = \frac{1}{bB^2}. \end{aligned}$$

Applying now Proposition 6.2 we get

$$\begin{aligned} & \mathbf{P}(|\xi_h| > \lambda\sigma_h(x), \mathcal{A}_h) \\ & \leq 4e \log(4B^3) \left(1 + 4r \sqrt{\frac{1+r}{1-\rho}} \lambda^2\right) \lambda e^{-\frac{\lambda^2}{2}}. \end{aligned} \quad (5.7)$$

5.4. AN UPPER BOUNDS FOR  $\zeta_h(1)$  AND  $\zeta_h(2)$ 

Note that both  $\zeta_h^{(1)}$ ,  $\zeta_h^{(2)}$  are linear combinations of elements of the form  $v_h \int_0^T \Psi(X_s^\varepsilon)[a(Y_s^\varepsilon) - \bar{a}] ds$ , where

- $v_h$  is any of  $v_{1,h}$ ,  $v_{2,h}$ ;
- $\Psi(X_s^\varepsilon)$  is any of  $\frac{(X_s^\varepsilon - x)^k}{h^k} K\left(\frac{X_s^\varepsilon - x}{h}\right)$ ,  $k = 0, 1, 2$ ;
- $a(Y_s^\varepsilon)$  is any of  $f(x, Y_s^\varepsilon)$ ,  $f_x(x, Y_s^\varepsilon)$ , and  $\bar{a} = \int a(y) p(y) dy$ , with  $p(\cdot)$  being the invariant density of the fast process.

Under the assumptions made, the function  $\Psi(u)$  is bounded by 1 and twice continuously differentiable: there exists a constant  $C_1$  such that

$$|\Psi(u)| \leq 1 \text{ and } |\dot{\Psi}(u)| + |\ddot{\Psi}(u)| \leq C_1 \quad \forall u.$$

Next, on the set  $\mathcal{A}_h$  it holds  $v_{1,h}^2 \leq \rho v_{0,h} v_{2,h} \leq \rho r v_{2,h}^2$  and  $v_{2,h} \leq bB(Th)^{-1}$ , so that, taking into account  $Th \geq 1$ , it suffices to bound only

$$U_T^\varepsilon = \int_0^T \Psi(X_s^\varepsilon)[a(Y_s^\varepsilon) - \bar{a}] ds.$$

We apply a large deviation type estimate for the two scaled diffusion model (1.1), (1.2) from Liptser and Spokoiny (1997) adapted to the case considered.

**PROPOSITION 5.1.** *Suppose  $(A_s)$  and  $(A_f)$ . If  $T = T_\varepsilon$  and  $\lim_{\varepsilon \rightarrow 0} T_\varepsilon \varepsilon = 0$ , then for every positive  $z > 0$  and  $0 < \kappa < 1/2$*

$$\lim_{\varepsilon \rightarrow 0} (\varepsilon T_\varepsilon)^{1-2\kappa} \log \mathbf{P} \left( (\varepsilon T_\varepsilon)^{-\kappa} |U_{T_\varepsilon}^\varepsilon| > z \right) \leq -\frac{z^2}{2\gamma},$$

where

$$\begin{aligned} \gamma &= \int_R \vartheta^2(y) G^2(y) p(y) dy, \\ \vartheta(y) &= \frac{2}{G^2(y) p(y)} \int_\infty^y [a(u) - \bar{a}] p(u) du. \end{aligned}$$

**COROLLARY 5.1.** *For  $\varepsilon$  small enough and  $\kappa_1 < 1 - 2\kappa$*

$$\mathbf{P} (|U_{T_\varepsilon}^\varepsilon| > (\varepsilon T_\varepsilon)^\kappa) < e^{-(\varepsilon T_\varepsilon)^{-\kappa_1}}.$$

Applying this corollary with  $\kappa < 1/2$  and  $\kappa_1 < 1 - 2\kappa$ , we obtain for  $\varepsilon T$  small enough

$$\mathbf{P} \left( |\zeta_h^{(i)}| > 2(\varepsilon T)^\kappa \right) < 2 \exp \left( -\frac{1}{(\varepsilon T)^{\kappa_1}} \right), \quad i = 1, 2. \quad (5.8)$$

### 5.5. PROOF OF THEOREM 3.1

Summing up the decomposition (5.3) and the bounds (5.6), (5.7), (5.8), we get

$$\begin{aligned} \mathbf{P} \left( \left| \tilde{f}_h(x) - \bar{f}(x) \right| > c\Delta_h(x) + \lambda\sigma_h(x) + 2(\varepsilon T^\varepsilon)^\kappa, \mathcal{A}_h \right) \\ \leq 4e \log(4B^3) \left( 1 + 4r \sqrt{\frac{1+r}{1-\rho}} \lambda^2 \right) \lambda e^{-\lambda^2/2} + 4e^{-(\varepsilon T)^{-\kappa_1}}. \end{aligned}$$

This leads to the required bound from Theorem 3.1 for sufficiently small  $\varepsilon T$ .

### 5.6. PROOF OF THEOREM 4.1

Let  $h_{\text{id}}$  be a “good” bandwidth. We intend to show that

$$\begin{aligned} \left\{ \left| \hat{f}(x) - \bar{f}(x) \right| > (2\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x) \right\} \\ \subseteq \bigcup_{\eta \in \mathcal{H}(h_{\text{id}})} \left\{ \left| \tilde{f}_\eta(x) - \bar{f}(x) \right| > \lambda_1\sigma_\eta(x) \right\} \end{aligned}$$

where  $\mathcal{H}(h) = \{\eta \in \mathcal{H} : \eta \leq h\}$ . This statement is equivalent to saying that the inequality  $\left| \hat{f}(x) - \bar{f}(x) \right| > (2\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x)$  is impossible if

$$\left| \tilde{f}_\eta(x) - \bar{f}(x) \right| \leq \lambda_1\sigma_\eta(x), \quad \forall \eta \in \mathcal{H}(h_{\text{id}}). \quad (5.9)$$

Obviously

$$\begin{aligned} \left\{ \left| \hat{f}(x) - \bar{f}(x) \right| > (2\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x) \right\} \\ \subseteq \left\{ \left| \hat{f}(x) - \bar{f}(x) \right| > (2\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x), \hat{h} > h_{\text{id}} \right\} \\ + \{h_{\text{id}} \text{ is rejected}\}. \end{aligned}$$

We consider separately each event in the right side of this inequality.

It holds on the event  $\{\hat{h} > h_{\text{id}}\}$  in view of the definition of  $\hat{h}$

$$\left| \tilde{f}_{\hat{h}}(x) - \tilde{f}_{h_{\text{id}}}(x) \right| \leq \lambda_1\sigma_{h_{\text{id}}}(x) + \lambda_2\sigma_{\hat{h}}(x) \leq (\lambda_1 + \lambda_2)\sigma_{h_{\text{id}}}(x).$$

Next, by (5.9)

$$|\tilde{f}_{h_{\text{id}}}(x) - \bar{f}(x)| \leq \lambda_1 \sigma_{h_{\text{id}}}(x).$$

Hence,  $\{\hat{h} > h_{\text{id}}\}$  and (5.9) imply

$$\begin{aligned} \left| \hat{f}(x) - \bar{f}(x) \right| &\leq |\tilde{f}_{\hat{h}}(x) - \tilde{f}_{h_{\text{id}}}(x)| + |\tilde{f}_{h_{\text{id}}}(x) - \bar{f}(x)| \\ &\leq (2\lambda_1 + \lambda_2) \sigma_{h_{\text{id}}}(x). \end{aligned}$$

Now we study the event  $\{h_{\text{id}} \text{ is rejected}\}$ . By definition

$$\begin{aligned} \{h_{\text{id}} \text{ is rejected}\} &= \bigcup_{h \in \mathcal{H}(h_{\text{id}})} \bigcup_{\eta \in \mathcal{H}(h)} \left\{ |\tilde{f}_h(x) - \tilde{f}_\eta(x)| > \lambda_2 \sigma_h(x) + \lambda_1 \sigma_\eta(x) \right\}. \end{aligned}$$

Condition (5.9) yields for every pair  $\eta < h \in \mathcal{H}(h_{\text{id}})$

$$|\tilde{f}_h(x) - \tilde{f}_\eta(x)| \leq |\tilde{f}_h(x) - \bar{f}(x)| + |\tilde{f}_\eta(x) - \bar{f}(x)| \leq \lambda_1 (\sigma_h(x) + \sigma_\eta(x))$$

so that the event  $\{h_{\text{id}} \text{ is rejected}\}$  is impossible under (5.9) in view of  $\lambda_2 \geq \lambda_1$ .

It remains to bound the probability of the event in (5.9). With  $\lambda_\eta = \lambda_1 - c\Delta_\eta(x)/\sigma_\eta(x)$ , it holds by Theorem 3.1

$$\begin{aligned} &\mathbf{P} \left( |\tilde{f}_\eta(x) - \bar{f}(x)| > \lambda_1 \sigma_\eta(x) \right) \\ &= \mathbf{P} \left( |\tilde{f}_\eta(x) - \bar{f}(x)| > \lambda_\eta \sigma_\eta(x) + c\Delta_\eta(x) \right) \leq \alpha(\lambda_\eta) \end{aligned}$$

where  $\alpha(\lambda)$  is from (3.6) and hence,

$$\mathbf{P} \left( |\tilde{f}_\eta(x) - \bar{f}(x)| \leq \lambda_1 \sigma_\eta(x), \forall \eta \in \mathcal{H}(h_{\text{id}}) \right) \geq 1 - \sum_{\eta \in \mathcal{H}(h_{\text{id}})} \alpha(\lambda_\eta)$$

This completes the proof of the theorem.

## 6. Appendix. Deviation probabilities for martingales

In the Appendix we present two general results for continuous martingales. The first result describes some properties of real-valued martingales, while the second one deals with martingales valued in  $\mathbb{R}^2$ .

### 6.1. THE SCALAR CASE

Let  $M_t$  be a continuous martingale with  $M_0 = 0$  and with the predictable quadratic variation  $\langle M \rangle_t$ .

PROPOSITION 6.1. For every  $T > 0$ ,  $\vartheta > 0$ ,  $S \geq 1$  and  $\lambda \geq 1$

$$\mathbf{P} \left( |M_T| > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \leq 4\lambda\sqrt{e} (1 + \log S) e^{-\frac{\lambda^2}{2}}.$$

*Proof.* We use

$$\begin{aligned} & \mathbf{P} \left( |M_T| > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \\ & \leq \mathbf{P} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \\ & \quad + \mathbf{P} \left( M_T < -\lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right). \end{aligned}$$

We estimate separately each term in the right side of this inequality.

Given  $a > 1$ , introduce the geometric series  $\vartheta_k = \vartheta a^k$  and define the sequence of random events  $\mathcal{C}_k = \{\vartheta_k \leq \sqrt{\langle M \rangle_T} < \vartheta_{k+1}\}$ ,  $k = 0, 1, \dots$ . Then clearly

$$\begin{aligned} & \mathbf{P} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \\ & \leq \sum_{k \geq 0}^K \mathbf{P} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S, \mathcal{C}_k \right). \end{aligned} \quad (6.1)$$

where  $K$  is the integer part of  $\log_a S$ . We now bound each term in this sum. Let, with  $\gamma \in \mathbb{R}$ ,

$$Z_t(\gamma) = \exp \left( \gamma M_t - \frac{\gamma^2}{2} \langle M \rangle_t \right).$$

The random process  $Z_t(\gamma)$  is the continuous local martingale and, being positive, it is the supermartingale (see Problem 1.4.4 in Liptser and Shiriyayev (1986)). Therefore for every  $T > 0$ ,

$$\mathbf{E} Z_T(\gamma) \leq 1. \quad (6.2)$$

For fixed  $k$ , we pick  $\gamma_k = \frac{\lambda}{\vartheta_k}$  and use (6.2) for the inequality

$$1 \geq \mathbf{E} Z_T(\gamma_k) \mathbf{I} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \mathcal{C}_k \right)$$

which implies

$$\begin{aligned} 1 & \geq \mathbf{E} \exp \left( \frac{\lambda}{\vartheta_k} M_T - \frac{\lambda^2}{2\vartheta_k^2} \langle M \rangle_T \right) \mathbf{I} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \mathcal{C}_k \right) \\ & \geq \mathbf{E} \exp \left( \frac{\lambda^2}{\vartheta_k} \sqrt{\langle M \rangle_T} - \frac{\lambda^2}{2\vartheta_k^2} \langle M \rangle_T \right) \mathbf{I} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \mathcal{C}_k \right) \\ & \geq \mathbf{E} \exp \left\{ \inf_{\vartheta_k \leq v \leq \vartheta_{k+1}} \left( \frac{\lambda^2 v}{\vartheta_k} - \frac{\lambda^2 v^2}{2\vartheta_k^2} \right) \right\} \mathbf{I} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \mathcal{C}_k \right). \end{aligned}$$

It is easy to check that “ $\inf_{\vartheta_k \leq v \leq \vartheta_{k+1}}$ ” is attained at the point  $v = \vartheta_{k+1} = a\vartheta_k$  so that

$$\mathbf{P} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \mathcal{C}_k \right) \leq \exp \left\{ -\lambda^2 \left( a - \frac{a^2}{2} \right) \right\}.$$

Combining this bound with (6.1) and using  $K \leq \log_a S$ , we obtain

$$\begin{aligned} \mathbf{P} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \\ \leq (1 + \log_a S) \exp \left\{ -\lambda^2 \left( a - \frac{a^2}{2} \right) \right\}. \end{aligned}$$

Since the left hand side of this inequality does not depend on  $a$ , we may optimize the choice of  $a$  to minimize its right side. This leads to  $a = 1 + 1/\lambda$ . Then

$$\lambda^2 \left( a - \frac{a^2}{2} \right) = \lambda^2 \left\{ 1 + \frac{1}{\lambda} - \frac{1}{2} \left( 1 + \frac{1}{\lambda} \right)^2 \right\} = \frac{1}{2} (\lambda^2 - 1)$$

and, since  $\log(1 + 1/\lambda) \geq 1/(2\lambda)$  for  $\lambda \geq 1$ , we have  $\log_a S \leq 2\lambda \log S$ . Hence

$$\mathbf{P} \left( M_T > \lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \leq 2\sqrt{e}\lambda (1 + \log S) e^{-\frac{\lambda^2}{2}}.$$

In the similar way we obtain

$$\mathbf{P} \left( M_T < -\lambda \sqrt{\langle M \rangle_T}, \vartheta \leq \sqrt{\langle M \rangle_T} \leq \vartheta S \right) \leq 2\sqrt{e}\lambda (1 + \log S) e^{-\frac{\lambda^2}{2}}$$

and the assertion follows.

## 6.2. THE VECTOR CASE

Here, we consider continuous vector martingale  $M_t$  valued in  $\mathbb{R}^2$  with components  $M_{0,t}$  and  $M_{1,t}$ . We denote

$$V_{0,t} = \langle M_0 \rangle_t, \quad V_{1,t} = \langle M_0, M_1 \rangle_t, \quad V_{2,t} = \langle M_1 \rangle_t.$$

Let  $u$  be a random variable and

$$\sigma_t^2 = V_{0,t} - 2uV_{1,t} + u^2V_{2,t}.$$



For a fixed time moment  $T$  and constants  $\vartheta > 0$ ,  $S \geq 1$ ,  $\beta \geq 0$  and  $\rho \in (0, 1)$ , introduce the event

$$\mathcal{A}_T = \left\{ \begin{array}{l} \vartheta \leq \sigma_T^2 \leq \vartheta S \\ V_{1,T}^2 \leq \rho V_{0,T} V_{2,T} \\ |u| \leq \beta \end{array} \right\}. \quad (6.3)$$

**PROPOSITION 6.2.** *Let  $M_t$  be a martingale with values in  $\mathbb{R}^2$  such that  $V_{0,T} \geq V_{2,T}$ . Then, with  $\mathcal{A}_T$  from (6.3), it holds for every  $\lambda \geq \sqrt{2}$ ,*

$$\begin{aligned} & \mathbf{P}(|M_{0,T} - uM_{1,T}| > \lambda\sigma_T, \mathcal{A}_T) \\ & \leq 4e \log(4S) \left( 1 + 4\beta \sqrt{\frac{1+\beta}{1-\rho}} \lambda^2 \right) \lambda e^{-\frac{\lambda^2}{2}}. \end{aligned}$$

*Proof.* For fixed  $\beta$ ,  $\rho$ , and  $\lambda$  define  $\delta$  such that

$$\frac{2\delta(1+\beta)}{1-\rho} = \lambda^{-2} \quad (6.4)$$

and denote by  $D_\delta = \{\alpha_k = k\delta : k \in \mathbb{N}, |\alpha| \leq \beta\}$  the discrete grid with the step  $\delta$  in the interval  $[-\beta, \beta]$ .

Let  $\nu_+$  (respectively  $\nu_-$ ) be random variable valued in  $D_\delta$  which is closest to  $u$  from above (respectively from below). Then clearly

$$|\nu_\pm - u| \leq \delta. \quad (6.5)$$

$$|M_{0,T} - uM_{1,T}| \leq \max\{|M_{0,T} - \nu_- M_{1,T}|, |M_{0,T} - \nu_+ M_{1,T}|\}. \quad (6.6)$$

Let now  $\nu$  be one of  $\nu_-$  and  $\nu_+$ . Then by the construction  $|\nu - u| \leq \delta$ . Next we show that on the set  $\mathcal{A}_T$  it holds

$$1 - \lambda^{-2} \leq \frac{V_{0,T} - 2\nu V_{1,T} + \nu^2 V_{2,T}}{\sigma_T^2} \leq 1 + \lambda^{-2} \quad (6.7)$$

Indeed

$$\begin{aligned} \sigma_T^2 &= V_{0,T} - 2uV_{1,T} + u^2 V_{2,T} = V_{0,T} - \frac{V_{1,T}^2}{V_{2,T}} + V_{2,T} \left( u - \frac{V_{1,T}}{V_{2,T}} \right)^2 \\ &\geq \frac{V_{0,T}V_{2,T} - V_{1,T}^2}{V_{2,T}} \geq (1-\rho)V_{0,T} \end{aligned}$$

and using  $V_{2,T} \leq V_{0,T}$ , we get

$$\begin{aligned} \frac{|V_{1,T}|}{\sigma_T^2} &\leq \frac{\sqrt{\rho V_{0,T} V_{2,T}}}{(1-\rho)V_{0,T}} \leq \frac{\sqrt{\rho}}{1-\rho} \leq (1-\rho)^{-1}, \\ \frac{V_{2,T}}{\sigma_T^2} &\leq \frac{V_{2,T}}{(1-\rho)V_{0,T}} \leq (1-\rho)^{-1}. \end{aligned}$$

Since on the set  $\mathcal{A}$  it holds  $|u| \leq \beta$  and by construction  $\nu \leq \beta$  we obtain, using the definition (6.4) of  $\delta$ ,

$$\begin{aligned} & |V_{0,T} - 2uV_{1,T} + u^2V_{2,T} - (V_{0,T} - 2\nu V_{1,T} + \nu^2V_{2,T})| \\ & \leq 2|V_{1,T}||u - \nu| + V_{2,T}|u^2 - \nu^2| \\ & \leq 2\delta(1 - \rho)^{-1}\sigma_T^2 + 2\beta\delta(1 - \rho)^{-1}\sigma_T^2 \\ & = \sigma_T^2\lambda^{-2} \end{aligned}$$

and (6.7) follows.

Since on the set  $\mathcal{A}_T$  the value  $\sigma_T^2$  is between  $\vartheta$  and  $\vartheta S$ , we also get for  $\nu = \nu_{\pm}$

$$(1 - \lambda^{-2})\vartheta \leq V_{0,T} - 2\nu V_{1,T} + \nu^2V_{2,T} \leq (1 + \lambda^{-2})\vartheta S. \quad (6.8)$$

We now derive from (6.6), (6.7) and (6.8)

$$\begin{aligned} & \{M_{0,T} - uM_{1,T} > \lambda\sigma_T, \mathcal{A}_T\} \\ & \subseteq \left\{ M_{0,T} - \nu_- M_{1,T} > \frac{\lambda}{\sqrt{1 + \lambda^2}} \sqrt{V_T(\nu_-)}, \mathcal{A}_T \right\} \\ & \quad \cup \left\{ M_{0,T} - \nu_+ M_{1,T} > \frac{\lambda}{\sqrt{1 + \lambda^2}} \sqrt{V_T(\nu_+)}, \mathcal{A}_T \right\} \\ & \subseteq \bigcup_{\alpha \in D_\delta} \left\{ |M_{0,T} - \alpha M_{1,T}| > \frac{\lambda}{\sqrt{1 + \lambda^2}} \sqrt{V_T(\alpha)}, \mathcal{A}_{\alpha,T} \right\}, \end{aligned}$$

where

$$\begin{aligned} V_T(\alpha) &= V_{0,T} - 2\alpha V_{1,T} + \alpha^2 V_{2,T}, \\ \mathcal{A}_{\alpha,T} &= \{(1 - \lambda^{-2})\vartheta \leq V_T(\alpha) \leq (1 + \lambda^{-2})\vartheta S\}. \end{aligned}$$

Now, for every  $\alpha \in D_\delta$ , the process  $M_{0,t} - \alpha M_{1,t}$  is the continuous local martingale with  $\langle M_0 - \alpha M_1 \rangle_T = V_{0,T} - 2\alpha V_{1,T} + \alpha^2 V_{2,T}$ . Applying Proposition 6.1 and using the inequalities  $\lambda^2 \geq 2$  and  $\frac{\lambda^2}{1 + \lambda^{-2}} \geq \lambda^2(1 - \lambda^{-2}) = \lambda^2 - 1$ , we obtain

$$\begin{aligned} & \mathbf{P} \left( |M_{0,T} - \alpha M_{1,T}| > \frac{\lambda}{\sqrt{1 + \lambda^2}} \sqrt{V_T(\alpha)}, \mathcal{A}_{\alpha,T} \right) \\ & \leq 4 \frac{\lambda}{\sqrt{1 + \lambda^{-2}}} \left( 1 + \log \frac{(1 + \lambda^{-2})\vartheta S}{(1 - \lambda^{-2})\vartheta} \right) \exp \left( -\frac{\lambda^2}{2(1 + \lambda^{-2})} + \frac{1}{2} \right) \\ & \leq 4\lambda \left( 1 + \log \frac{3S}{2} \right) \exp \left( -\frac{\lambda^2}{2} + 1 \right). \end{aligned}$$

Since the number of different elements in  $D_\delta$  is at most  $1 + 2\beta\delta^{-1}$  and since  $\delta$  from (6.4) fulfills  $\delta^{-1} = \frac{2\lambda^2(1+\beta)}{1-\rho}$  we get

$$\mathbf{P}(|M_{0,T} - uM_{1,T}| > \lambda\sigma_T, \mathcal{A}_T)$$

$$\begin{aligned} &\leq 4e \left( 1 + \log \frac{3S}{2} \right) (1 + 2\beta\delta^{-1}) \lambda e^{-\frac{\lambda^2}{2}} \\ &\leq 4e \log(4S) \left( 1 + 4\beta \sqrt{\frac{1+\beta}{1-\rho}} \lambda^2 \right) \lambda e^{-\frac{\lambda^2}{2}} \end{aligned}$$

as required.

## References

1. Bhattacharya, R. N.; Goetze, F. (1995). Time-scales for Gaussian approximation and its breakdown under a hierarchy of periodic spatial heterogeneities. *Bernoulli* **1**, no.1-2, 81-123.
2. Brown, L.D. and Low, M.G. (1996). A constrained risk inequality with applications to nonparametric functional estimation. *Ann.-Statist.* **24** (1996), no. 6, 2524–2535.
3. Carrol, R.J. and Hall, P. (1988). Optimal rate of convergence for deconvoluting a density. *J. Amer. Statist. Assoc.* **83**, 1184–1186.
4. G. Collomb and P. Doukhan (1983). Estimation non parametrique de la fonction d'autoregression d'un processus stationnaire et phi melangeant: risques quadratiques pour la methode du noyau, *C. R. Acad. Sci., Paris, Ser. I*, **296**, 859-862 .
5. Delyon, B. and Juditsky, A. (1997). On minimax prediction for nonparametric autoregressive models. Unpublished manuscript.
6. P. Doukhan and M. Ghindes (1980). Estimations dans le processus “ $X_{n+1} = f(X_n) + \epsilon_n$ ”, *C. R. Acad. Sci., Paris, Ser. A* **291**, 61-64.
7. Doukhan, P. and Tsybakov, A.B. (1993). Nonparametric recurrent estimation in nonlinear ARX models. *Problemy-Peredachi-Informatsii* **29**, no. 4, 24–34. Translation: *Problems Inform. Trans.* **29**, no. 4, 318–327.
8. Fan, J. and Truong, Y.K. (1993). Nonparametric regression with errors in variables. *Ann. Statist.* **21**, 1900–1925.
9. Fan, J. and Gijbels, I. (1996). *Local polynomial modelling and its applications*. Chapman and Hall, London.
10. Freidlin, M.I., Wentzell A.D. (1984). *Random Perturbations of Dynamical Systems*. Springer. N.Y
11. Genon-Catalot, V., Laredo, C. and Picard, D. (1992). Nonparametric estimation of the diffusion coefficient by wavelet methods, *Scand. J. Statist.* **19**, 317-335.
12. W. Härdle and P. Vieu (1992). Kernel regression smoothing of time series”, *J. Time Ser. Anal.* **13**, No.3, 209-232.
13. Grama, I. and Nussbaum, M. (1998). Asymptotic equivalence for nonparametric generalized linear models, *Prob. Theory and Rel. Fields*, **111**, 167–214.
14. Ibragimov, I.A. and Khasminskii, R.Z. (1981). *Statistical Estimation: Asymptotic Theory* Springer, New York.

15. Katkovnik, V. Ja. (1985). *Nonparametric Identification and Data Smoothing: Local Approximation Approach*. Nauka, Moscow (in Russian).
16. Khasminskii, R.Z. (1966). On stochastic processes defined by differential equations with small parameter. *Theory Probab. Appl.* **11**, 219–228.
17. H.J.Kushner (1990). *Weak Convergence Methods and Singularly Perturbed Stochastic Control and Filtering Problems*. Birkhäuser 1990.
18. Kutoyants, Yu.A. (1984a). On nonparametric estimation of trend coefficients in a diffusion process. Collection: *Statistics and control of stochastic processes*, Moscow, 230–250.
19. Kutoyants, Yu.A. (1984b). Parameter estimation for stochastic processes. Translated from the Russian and edited by B. L. S. Prakasa Rao. *R & E Research and Exposition in Mathematics*, **6**. Heldermann Verlag, Berlin.
20. Kutoyants, Yu.A. (1994). *Identification of dynamical systems with small noise*. Kluwer, Dordrecht.
21. Lepski, O. (1990). One problem of adaptive estimation in Gaussian white noise. *Theory Probab. Appl.* **35**, no. 3, 459–470.
22. Lepski, O. and Levit, B. (1997). Efficient adaptive estimation of infinitely differentiable function. *Math. Methods of Statistics*, submitted.
23. Lepski, O., Mammen, E. and Spokoiny, V. (1997). Ideal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selection. *Annals of Statistics*, **25**, no.3, 929–947.
24. Lepski, O. and Spokoiny, V. (1997). Optimal pointwise adaptive methods in nonparametric estimation. *Annals of Statistics*, **25**, no.6,
25. Liptser, R. and Shiryaev, A. (1989). *Theory of Martingales*. Kluwer Acad. Publ. 1989.
26. R.S. Liptser, W.J. Runggaldier and M. Taksar (1996). Deterministic approximation for stochastic control problem. *SIAM J. on Control and Optimization*. **34** 161–178.
27. Liptser, R. and Spokoiny, V. (1999). Moderate deviations type evaluation for integral functionals of diffusion processes. *Electronic J. of Probability* **4** no. 17, 1–25.
28. Mercurio, D. and Spokoiny, V. (2000). Statistical inference for time-inhomogeneous volatility models. Unpublished manuscript.
29. Spokoiny, V. (1996). Adaptive hypothesis testing using wavelets. *Annals of Stat.*, **24**, no. 6, 2477–2498.
30. Milstein, G. and Nussbaum, M. (1994). Nonparametric estimation of a nonparametric diffusion model, *Prob. Theory and Rel. Fields*. To appear.
31. Tsybakov, A. (1986). Robust reconstruction of functions by the local approximation. *Prob. Inf. Transm.*, **22**, 133-146.
32. Veretennikov, A. Yu. (1991) On the averaging principle for systems of stochastic differential equations. *Math. USSR Sborn.*, **69**, No. 1, 271-284.
33. Veretennikov, A. Yu. (1992) On large deviations for ergodic empirical measures. *Topics in Nonparametric Estimation. Advances in Soviet Mathematics*, AMS **12**, 125-133.