

Hong, Seok Young; Linton, Oliver Bruce

**Working Paper**

## Asymptotic properties of a Nadaraya-Watson type estimator for regression functions of infinite order

cemmap working paper, No. CWP53/16

**Provided in Cooperation with:**

Institute for Fiscal Studies (IFS), London

*Suggested Citation:* Hong, Seok Young; Linton, Oliver Bruce (2016) : Asymptotic properties of a Nadaraya-Watson type estimator for regression functions of infinite order, cemmap working paper, No. CWP53/16, Centre for Microdata Methods and Practice (cemmap), London, <https://doi.org/10.1920/wp.cem.2016.5316>

This Version is available at:

<https://hdl.handle.net/10419/189788>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

# Asymptotic properties of a Nadaraya-Watson type estimator for regression functions of infinite order

---

Seok Young Hong  
Oliver Linton

The Institute for Fiscal Studies  
Department of Economics, UCL

**cemmap** working paper CWP53/16

# Asymptotic properties of a Nadaraya-Watson type estimator for regression functions of infinite order\*

Seok Young Hong<sup>†</sup>      Oliver Linton<sup>‡</sup>

*University of Cambridge*

November 23, 2016

## Abstract

We consider a class of nonparametric time series regression models in which the regressor takes values in a sequence space and the data are stationary and weakly dependent. We propose an infinite dimensional Nadaraya-Watson type estimator with a bandwidth sequence that shrinks the effects of long lags. We investigate its asymptotic properties in detail under both static and dynamic regressions contexts. First we show pointwise consistency of the estimator under a set of mild regularity conditions. We establish a CLT for the estimator at a point under stronger conditions as well as for a feasibly studentized version of the estimator, thereby allowing pointwise inference to be conducted. We establish the uniform consistency over a compact set of logarithmically increasing dimension. We specify the explicit rates of convergence in terms of the Lambert W function, and show that the optimal rate that balances the squared bias and variance is of logarithmic order, the precise rate depending on the smoothness of the regression function and the dependence of the data in a non-trivial way.

KEYWORDS: Functional Regression; Nadaraya-Watson estimator; Curse of infinite dimensionality; Near Epoch Dependence.

---

\*Part of this paper was written while the first author visited Institut de Mathématiques, Toulouse in November 2013. We thank Philippe Vieu for his kind hospitality and many helpful advice throughout writing this paper. We also thank John Aston, Xiaohong Chen, Jeroen Dalderop, Paul Doukhan, Jiti Gao, Hayden Lee, Mikhail Lifshits, Richard Nickl, Alexei Onatski, and Peter Phillips for their helpful comments. Special thanks to Alexey Rudenko for providing an original Russian photocopy of Sytaya (1974), and to Hyungjin Lee for translating the paper. Financial support from the European Research Council (ERC-2008AdG-NAMSEF) is gratefully acknowledged.

<sup>†</sup>Statistical Laboratory, Department of Pure Mathematics and Mathematical Statistics, Faculty of Mathematics, University of Cambridge, United Kingdom, email: syh30@cam.ac.uk

<sup>‡</sup>Faculty of Economics, University of Cambridge, United Kingdom, e-mail: obl20@cam.ac.uk

# 1 Introduction

Nonparametric modelling is a well established practical tool for analyzing time series data; see for example Härdle (1990), Bosq (1996), or Fan and Yao (2003) for a comprehensive review. A major advantage of this approach is that the relationship between the explanatory variables under study, denoted by  $X = (X_1, \dots, X_d)^\top$ , and the response, say  $Y$ , can be modelled without assuming any restrictive parametric or linear structures. One issue with allowing for this extended flexibility is the *curse of dimensionality*; Stone (1980, 1982) showed that given a fixed measure of smoothness  $\beta$  allowed on the regression function the best achievable convergence rate (in minimax sense)  $n^{-\beta/(2\beta+d)}$  deteriorates dramatically as the dimension/order  $d$  increases.

In a time series context it is often reasonable to model the dependence upon the infinite past. For example, the  $\text{AR}(d)$  and  $\text{ARX}(d)$  models with  $d = \infty$  naturally extends those classical linear models, and enables the influence of *all past information* to be taken into account, thereby allowing for maximal flexibility with regard to the dynamic structure. It can also be very useful for several semiparametric applications, and for testing the martingale hypothesis or the efficient market hypothesis in economics, where the conditional mean given all past information  $E(Y_t|\mathcal{F}_{t-1})$  is the object of main interest. It is desirable to be able to nest the linear  $\text{AR}(\infty)$  models, and this is one of our aims. Not restricting the number of conditioning variables also has an advantage of avoiding the econometrician's *a priori* choice of the order  $d$  based on some order determination principles whose validity is often subject to question in practical situations. For these reasons, we study a class of nonparametric time series regression models of infinite order that covers both static and dynamic regression cases, and in particular includes the autoregression framework as a special case.

Pagan and Hong (1990) proposed studying the nonparametric regression case where  $d \rightarrow \infty$  in the context of the econometric analysis of risk models. Linton and Sancetta (2009) tackled the estimation problem in the context of an autoregressive model and established uniform almost sure consistency for stationary ergodic sample observations. There is a vast literature on statistical research on functional data (typical examples include curves and images), which are infinite-dimensional in nature. Ferraty and Vieu (2002) first studied the case where the regressor was function valued. Masry (2005) provided a rigorous treatment of nonparametric regression with dependent functional data in which  $X$  lies in a general semi-metric space, establishing the central limit theorem. Mas (2012) derived the minimax rate of convergence for nonparametric estimation of the regression function with strictly independent and identically distributed covariates. Ferraty and Vieu (2006) detailed a number of extensions and gave an overview of nonparametric approaches in the functional statistics literature. Geenens (2011) gave an up-to-date accessible summary of the literature on nonparametric functional regression, and introduced the term *curse of infinite dimensionality*, which reflects the evident difficulties in nonparametric estimation of infinite-dimensional objects due to extreme data sparsity. We discuss in the next section the difference between the functional data framework and our discrete time framework.

One major challenge in the infinite-dimensional setting is that the usual notion of density  $p(\cdot)$  does not exist. Since there is no  $\sigma$ -finite Lebesgue measure in infinite-dimensional spaces, the Lebesgue density (with respect to the infinite product of probability measures) of the regressor cannot be defined via the Radon-Nikodym theorem. Consequently, standard asymptotic arguments for kernel estimators are no longer valid, for example, Bochner's lemma whereby under suitable regularity conditions, for  $j = 1, 2$

$$\frac{1}{h^d} E \left[ \mathcal{K}^j \left( \frac{x - X}{h} \right) \right] = \int \mathcal{K}^j(u) p(x - uh) du \rightarrow p(x) \|\mathcal{K}\|_j^j \quad \text{as } h \rightarrow 0 \quad (1)$$

where  $\mathcal{K}$  is a multivariate kernel function (see subsection 2.2 below). Hence, classical limiting theories in nonparametric literature cannot be readily extended.

In this paper, we consider an infinite dimensional analogue of the classical nonparametric regression approach. We propose a Nadaraya-Watson type kernel estimator and investigate its large sample properties. In particular, we establish both pointwise and uniform consistency of the estimator and establish its asymptotic normality under both static and dynamic regression contexts under  $\alpha$ -mixing and near epoch dependent sample observations. We impose some regularity conditions on the vector bandwidth sequence, and we derive the rate of convergence via specifying the small deviation probabilities. The pointwise and uniform rates are logarithmic, which reflects the difficulty of capturing nonparametrically the effect of an infinite number of lags. Our pointwise rate is consistent with the rate in Mas (2012) who derived under strict cross-sectional and temporal independence.

For notations, we define  $a_n \simeq b_n$  by  $a_n = b_n + o(1)$ , and  $c_n \sim d_n$  by equivalence of order between the two sequences  $c_n$  and  $d_n$ . Also,  $f \preceq g$  means there exists some constant  $c > 0$  such that  $\lim_{n \rightarrow \infty} f(n)/g(n) \leq c$ . The term 'stationarity' is taken to mean strict stationarity. Throughout,  $C$  (or  $C'$ ,  $C''$ ) refers to some generic constant that may take different values in different places unless defined specifically otherwise.

## 2 Some Preliminaries

We consider the regression model

$$Y = m(X) + \varepsilon, \quad (2)$$

where the regressor  $X = (X_1, X_2, \dots)^\top$  is a random element taking values in some sequence space  $S$ , the response  $Y$  is a real-valued variable, and the stochastic error  $\varepsilon$  is such that  $E(\varepsilon|X) = 0$  a.s. The objective is to estimate the Borel function  $m(\cdot) = E(Y|X = \cdot)$  based on  $n$  random samples observed from a strictly stationary data generating process  $\{(Y_t, X_t) \in \mathbb{R} \times S\}_{t \in \mathbb{Z}}$  having some weak dependence structure (see section 2.1 below).

This setting is related to the usual framework adopted for functional data, which has been widely studied by statisticians, see Ramsey and Silverman (2002). Recently, successful attempts have been made to develop theories for nonparametric inference in the functional statistics literature; Ferraty

and Romain (2010) gives a comprehensive review. A major issue in this field of research lies in extending the statistical theories applicable to  $\mathbb{R}^d$  to function spaces. In this literature, attention is usually on smooth functions that are approximated and reconstructed from finely discretised grids on some compact interval. In contrast, the setup in our model (2) can be viewed as looking at a countable number of discrete observations. Such a difference is reflected by the fact that the observed data is taken to be a discrete process  $X = (X_s)$  with unbounded  $s \in \mathbb{Z}^+$  so that  $S = \{f|f : \mathbb{N} \rightarrow \mathbb{R}\}$ , rather than  $X = (X(s))$  with  $s \in [0, T]^k$  so that  $S = \{f|f : [0, T]^k \subset \mathbb{R}^k \rightarrow \mathbb{R}\}$ , e.g. curves if  $k = 1$ , images if  $k \geq 2$ . The discrete nature of our setting has several fundamental distinctive features that allow us to look further into many specific practical issues.

An immediate consequence of our framework is that the tuning parameter can be imposed on each and every dimension, allowing one to control the marginal influence of the regressors. For instance when it is sensible to postulate that the influence of distant covariates is getting monotonically downweighted, one may set the marginal bandwidths to increase in the lag horizon so as to impose higher amount of smoothing. Depending on the nature of the regressor,  $S$  may be taken as the space of all infinite real sequences  $\mathbb{R}^\infty := \prod_{j=1}^\infty \mathbb{R}_j$  formed by taking Cartesian products of the reals, or its various linear subspaces such as  $\ell_\infty, \ell_p, c$ . We propose to take  $S = \mathbb{R}^\infty$  so as to refrain from imposing any prior restrictions with regard to the choice of the regressor; for example, taking  $S$  to be the space of bounded sequences excludes the possibility of the regressors with infinite supports (e.g. Gaussian process).

## 2.1 Dependence structure and leading examples

A distinctive characteristic of time series data is temporal dependence between observations. In the nonparametric time series literature, Rosenblatt (1956)'s  $\alpha$ -mixing has been the *de facto* standard choice due to it being the weakest among the class of mixing-type asymptotic independence conditions. Roussas (1990) established pointwise and uniform consistency of the local constant estimator under this condition, respectively, while Fan and Masry (1992) established asymptotic normality. The  $\alpha$ -mixing condition has also been widely used in the context of dependent functional observations, see for instance Ferraty et al. (2010), Masry (2005), and Delsol (2009).

**DEFINITION 1.** A stochastic process  $\{Z_t\}_{t=1}^\infty$  defined on some probability space  $(\Omega, \mathcal{F}, P)$  is called  $\alpha$ -mixing (cf. 'jointly'  $\alpha$ -mixing if  $Z_t$  is  $\mathbb{R}^d$ -valued, with  $d \in (1, \infty]$ ) if

$$\alpha(r) := \sup_{A \in \mathcal{F}_{-\infty}^t, B \in \mathcal{F}_{r+t}^\infty} |P(A \cap B) - P(A)P(B)|$$

is asymptotically zero as  $r \rightarrow \infty$ , where  $\mathcal{F}_a^b$  is the  $\sigma$ -algebra generated by  $\{Z_s; a \leq s \leq b\}$ . In particular, we say the process is algebraically (respectively exponentially)  $\alpha$ -mixing if there exists some  $c, k > 0$  such that  $\alpha(r) \leq cr^{-k}$  (respectively if there exists some  $\gamma, \varsigma > 0$  such that  $\alpha(r) \leq \exp(-\varsigma r^\gamma)$ ).

The popularity of the  $\alpha$ -mixing condition (note the modifier  $\alpha$ - will occasionally be omitted if no confusion is likely) in the literature stems from the fact that it is easy to work with, see e.g. Doukhan (1994) or Rio (2000) for a comprehensive survey. However, there are several limitations that have been pointed out in the literature. First, it is a rather strong technical condition that is hard to verify in practice. Second, some basic processes are not mixing. e.g. AR(1) with Bernoulli innovations, Andrews (1984).

We turn to our setting. In the static regression case it is appropriate to assume the mixing condition, but in the dynamic case this condition is not generally applicable as we now explain. Recall that the object of estimation is the conditional mean  $E(Y_t|\mathcal{F})$ , cf. (2), where the information set  $\mathcal{F}$  is determined by the nature of the conditioning variables. There are two leading cases: the first case is the static regression where the information set is taken to mean  $\sigma(X_{jt}; j = 1, 2, \dots)$ , the  $\sigma$ -algebra generated by the exogenous marginal regressors. The second case is the autoregression, where  $X_{tj} = Y_{t-j}$  for all  $j$ , in which case  $\mathcal{F} = \mathcal{F}_{t-1}$  represents  $\sigma(Y_s; s \leq t-1)$ , the  $\sigma$ -algebra generated by the sequence of the lags of the response  $(Y_s)_{s \leq t-1}$ . In fact, as for the latter framework we may consider a more general setup, i.e. a dynamic regression, where the information set is taken to be  $\mathcal{F} = \sigma(X_{js}, Y_s; s \leq t-1)$  for some  $j$ . Details are formally given in Assumptions A below.

In the static regression case the usual joint  $\alpha$ -mixing condition can be assumed on the sample data  $\{Y_t, X_t\}$  as is usually done; since marginal regressors are observed at the same time  $t$ :  $X_t = (X_{1t}, X_{2t}, \dots)^\top$ , assuming joint dependence does not require additional adjustments. Indeed, joint mixing implies that both marginal component processes and any measurable function thereof are mixing.<sup>1</sup> In this paper, we do not necessarily require independence between component processes  $\{X_{jt}\}$ ,  $j = 1, 2, \dots$ ; later we specify to what extent some dependence can be allowed (see Assumption C2). It will turn out that the requirement is mild and allows sufficient generality in application.

Moving on to the dynamic regression setting, since the regressors are taken to be the lags of the response and/or a covariate, measurable functions of  $X_t$  depend on infinite time-lags and hence are *not* necessarily mixing.<sup>2</sup> Therefore an alternative set of dependence conditions is necessary to establish asymptotic theories for the second framework. We shall adopt the notion of near epoch dependence due to Ibragimov (1962) as for the dynamic regression setting and deal with two leading cases separately.

**DEFINITION 2.** *A stochastic process  $\{Z_t\}_{t=1}^\infty$  defined on some probability space  $(\Omega, \mathcal{F}, P)$  is called near-epoch dependent or stable in  $L_2$  with respect to a strictly stationary  $\alpha$ -mixing process  $\{\eta_t\}$  if the stability coefficients  $v_2(r) := E|Z_t - Z_{t,(r)}|^2$  is asymptotically zero as  $r \rightarrow \infty$ , where  $Z_{t,(r)} = \Psi_r(\eta_t, \dots, \eta_{t-r+1})$  for some Borel function  $\Psi_r : \mathbb{R}^r \rightarrow \mathbb{R}$ .*

---

<sup>1</sup>The converse is not necessarily true unless the marginal processes are independent to each other, see Bradley (2005, Section 5).

<sup>2</sup>Except for some very special cases; Davidson (1994, Theorem 14.9) gives a set of technical conditions under which a process with infinite (linear) temporal dependence is  $\alpha$ -mixing.

A process that is *near epoch* dependent on a mixing sequence is influenced primarily by the “recent past” of the sequence and hence asymptotically resembles its dependence structure; see e.g. Billingsley (1968), Davidson (1994), or Lu (2001) for details. Andrews (1995) established uniform consistency of kernel regression estimators under near epoch dependence conditions. Following the usual convention, e.g. Bierens (1983), we shall take  $\Psi_r(\eta_t, \dots, \eta_{t-r+1}) \equiv E(Z_t | \eta_t, \dots, \eta_{t-r+1})$ . In section 3 it will be shown that under suitable conditions similar asymptotic theories can be derived for both static and dynamic regression frameworks.

## 2.2 Local Weighting

In this section we fix the notions of local weighting and the measure of closeness between the data objects. Let  $K : [0, \infty) \rightarrow [0, \infty) =: \mathbb{R}_+$  be a univariate density function and for an element  $u$  of a normed sequence space, let

$$\mathcal{K}(u) := K(\|u\|). \quad (3)$$

In our setting the properties of  $K$  are crucially important. We now group the kernel functions into three subcategories depending on how they are generated. The first two, referred to as Type-I and Type-II kernels in Ferraty and Vieu (2006) generalize the usual ‘window’ kernels and monotonically decreasing kernels in finite dimension, respectively. Both types of kernels are continuous on a compact support  $[0, \lambda]$ .

**DEFINITION 3.** *A function  $K : [0, \infty) \rightarrow [0, \infty)$  is called a kernel of type-I if it integrates to 1, and if there exist real constants  $C_1, C_2$  (with  $0 < C_1 < C_2$ ) for which*

$$C_1 1_{[0, \lambda]}(u) \leq K(u) \leq C_2 1_{[0, \lambda]}(u), \quad (4)$$

*where  $\lambda$  is some fixed positive real number. A function  $K : [0, \infty) \rightarrow [0, \infty)$  is called a kernel of type-II if it satisfies (4) with  $C_1 \equiv 0$ , and is continuous on  $[0, \lambda]$  and differentiable on  $(0, \lambda)$  with the derivative  $K'$  that satisfies*

$$C_3 \leq K'(u) \leq C_4$$

*for some real constants  $C_3, C_4$  such that  $-\infty < C_3 < C_4 < 0$ .*

The definition above suggests that the uniform kernel on  $[0, \lambda]$  is a type-I kernel, and the Epanechnikov, Biweight and Bartlett kernels belong to the class of Type-II kernels. Some of those with semi-infinite support, for example (one-sided) Gaussian, are covered by the last group, which we will call the Type-III kernels.

**DEFINITION 4.** *A function  $K : [0, \infty) \rightarrow [0, \infty)$  is a kernel of type-III if it integrates to 1, and if it is of exponential type; that is,  $K(r) \propto \exp(Cr^\beta)$  for some  $\beta$  and  $C$ .*



## 2.3 Small deviations

The *small ball (or small deviation) probability* plays a crucial role in establishing the asymptotic theory. Let  $S^*$  be a sequence space equipped with some norm  $\|\cdot\|$ ; then the small ball probability of an  $S^*$ -valued random element  $Z$  is a function defined as

$$\varphi_z(h) := P(\|z - Z\| \leq h), \quad (5)$$

where  $h \in \mathbb{R}_+$ . We shall call the probability *centered* if  $z = 0$  (in which case we write  $\varphi(h)$ ), and *shifted* (with respect to some fixed point  $z \in S^*$ ) if otherwise. The relation between the two quantities cannot be explicitly specified in general, and will be given in terms of the Radon-Nikodym derivative (See Assumption D1 below).

The name *small ball* stems from the fact that we are interested in the asymptotic behaviour of  $\varphi_z(h)$  as  $h$  tends to zero. The function can be thought of as a measure for how much the observations are densely *packed* or *concentrated* around the fixed point  $z$  with respect to the associated norm and the reference distance  $h$ . From the definition it is straightforward to see that  $\varphi_z(h) \rightarrow 0$  as  $h \rightarrow 0$ , and that  $n\varphi_z(h)$  is an approximate count of the number of observations whose influence is taken into account in the smoothing procedure. When  $Z$  is a continuous random vector of fixed dimension  $d$  with density  $p(\cdot) > 0$ , it can be readily shown that the shifted small ball probability (with respect to the usual Euclidean norm) is given by

$$\varphi_z(h) = V_d h^d p(z) = O(h^d), \quad (6)$$

where  $V_d = \pi^{d/2}/\Gamma(d/2 + 1)$  is the volume of the  $d$ -dimensional unit sphere.

However, when  $Z$  takes values in an infinite-dimensional normed space, it is difficult to specify the exact form of the small ball probability, and its behaviour varies depending heavily on the nature of the associated space and its topological structure. Due to the non-equivalence of norms in infinite dimensional spaces, it is intuitively clear that the “speed” at which  $\varphi_z(h)$  converges to zero is affected by the choice of the norm  $\|\cdot\|$ . Nonetheless, a rapid decay is expected in general irrespective of the choice of the norm due to the extreme sparsity of data in infinite-dimensional spaces.

One possible example of  $S^*$  is  $(\ell_r, \|\cdot\|_r)$ , the space of  $r$ -th power summable sequence equipped with the  $\ell_r$ -norm; the centred small ball behaviour of sums of weighted i.i.d. random variables is widely studied in the literature, see for example Borovkov and Ruzankin (2008) and references therein. In this paper, we will focus our main attention on the case of  $r = 2$  (and take  $\|\cdot\|$  to mean  $\|\cdot\|_2$  unless specified otherwise). Nevertheless, it is worth noting that the results derived in this paper can be extended to the case of  $r > 2$  as long as the regularity conditions are adjusted appropriately.

Writing the expected value of the kernel in terms of the small ball probability

$$EK\left(\frac{z - Z}{h}\right) = EK\left(\frac{\|z - Z\|}{h}\right) = \int K(u) dP_{\|z - Z\|/h}(u) = \int K(u) d\varphi_z(uh), \quad (7)$$

we are able to bypass the difficulties mentioned in the introduction, and to establish the convergence of the integrals without explicitly requiring the existence of the Lebesgue density.

LEMMA 1. Ferraty and Vieu (2006, Lemma 4.3 & 4.4). *Suppose  $\|\cdot\|$  is some semi-norm defined on a function space. If  $K$  is type-I, then it satisfies*

$$C_1^j \leq \frac{1}{\varphi_z(h\lambda)} \int_0^\lambda K^j(v) d\varphi_z(vh) \leq C_2^j, \quad j = 1, 2 \quad (8)$$

where  $C_1, C_2 > 0$  are as defined in Definition 3. When the kernel  $K$  is type-II, if

$$\exists \varepsilon_0 > 0, C_5 > 0 \text{ s.t. } \forall \varepsilon < \varepsilon_0, \int_0^\varepsilon \varphi_x(u) du > C_5 \varepsilon \varphi_x(\varepsilon) \quad (9)$$

then we have

$$C_6^j \leq \frac{1}{\varphi_z(h\lambda)} \int_0^\lambda K^j(v) d\varphi_z(vh) \leq C_7^j, \quad j = 1, 2 \quad (10)$$

where the constants  $C_6 = -C_5 C_4$  and  $C_7 = \sup_{s \in [0, \lambda]} K(s)$  are strictly positive.

Under the regularity conditions of Lemma 1, (8) and (10) hold for every  $h > 0$ , so it follows that for any kernels of type-I and II:

COROLLARY 1. *If the kernel  $K$  is either type-I or type-II, then for  $j = 1, 2$  we have*

$$\frac{1}{\varphi_z(h\lambda)} E \left[ \mathcal{K}^j \left( \frac{z - Z}{h} \right) \right] \longrightarrow \xi_j \quad \text{as } h \rightarrow 0^+, \quad (11)$$

where  $\xi_1$  and  $\xi_2$  are some strictly positive real constants.

This result can be seen as an infinite-dimensional analogue of Bochner's lemma (1): i.e., for  $Z \in \mathbb{R}^d$ ,  $h^{-d} E \mathcal{K}((z - Z)/h) \rightarrow p(z) > 0$ . It is obvious that  $\xi_j$  is bounded below and above by  $C_1^j$  and  $C_2^j$ , respectively (or  $C_6^j$  and  $C_7^j$  depending on the choice of the kernel). With specific choices of kernels and regressors we may be able to specify the exact values of the constants in some certain cases. For example, it is straightforward to see that  $\xi_1 = 1/\lambda$  and  $\xi_2 = 1/\lambda^2$  when  $K$  is uniform kernel.

REMARKS. (i) Lemma 1 reveals the importance of condition (9) in constructing the asymptotics when the kernel is of type-II. Whereas the condition is widely assumed in the functional statistics literature for that reason, Azais and Fort (2013) proved that it necessarily restricts the variable  $Z$  to be of finite dimension. In other words, whenever (9) is valid, the topology that governs the concentration properties of  $Z$  accounts effectively only for finite dimension. An example (cf. Section 13.3.3 of Ferraty and Vieu (2006)) includes the case where  $Z$  is associated with the semi-norm  $\|y\| := (y_1, \dots, y_p, 0, 0, \dots)$  for some positive integer  $p < \infty$ , where  $y \in \mathbb{R}^\infty$ . Since this severely

restricts the applicability of our paper, we shall not consider the case of Type-II kernels. (ii) A natural question one may then ask is whether (11) would hold for kernels with semi-infinite support such as the Type-III kernels. In the finite  $\mathbb{R}^d$ -framework, it is well known that a set of assumptions including  $\|u\|^d K(u) \rightarrow 0$  as  $u \rightarrow \infty$  is sufficient for showing (1), see for instance Parzen (1962, Theorem 1A) and Pagan and Ullah (1999, Lemma 1). However, in the infinite-dimensional setting the answer is negative in most usual cases where the kernel is of exponential type (e.g. Gaussian kernel). Whereas the lower bound of the limit can be easily constructed via Chebyshev's inequality: with reference to Definition 4, writing  $V = \|z - Z\|^\beta$ ,  $\delta = h^\beta$  and letting  $c_\delta$  be some function of  $\delta$  we have

$$(0 <) \exp(-c_\delta \delta) \leq [P(V \leq \delta)]^{-1} E \exp(-c_\delta V). \quad (12)$$

So the upper bound may not exist, and the rate at which the small ball probability decays to zero may dominate the speed at which the integral (7) converges to zero. This claim cannot be formally verified for all general cases because (as aforementioned) there is no unified result for the asymptotic behaviour of small deviations available. Nevertheless, the idea can be sketched in the common case where the asymptotics of the distribution function (i.e. small deviation) is of exponential order:  $P(V \leq \delta) \sim \exp(-C\delta^{-\theta})$  as  $\delta \rightarrow 0$  for some constants  $C$  and  $\theta > 0$ . By de Bruijn's exponential Tauberian theorem (see Bingham et al. (1987), Li (2012)), a necessary and sufficient condition for such a case is the following limiting behaviour of the Laplace transform near infinity:

$$E[\exp(-c_\delta V)] \sim \exp\left(-C' \cdot c_\delta^{\theta/(1+\theta)}\right) \quad \text{as } c_\delta \rightarrow \infty$$

for some constant  $C' > 0$ . With  $V = \|z - Z\|^2$ ,  $\delta = h^2$ ,  $c_\delta = 2^{-1}h^{-2}$  (which corresponds to the case of the Gaussian kernel) the difference in the order of convergence suggests that the right hand side of (12) is unbounded, and that the limit (11) diverges. Due to this reason, we shall confine our attention to compactly supported kernels in this paper.

## 2.4 Bandwidth Matrix and covariates

We will be estimating the regression operator at a point  $x \in \mathbb{R}^\infty$  with an  $\mathbb{R}^\infty$ -dimensional regressor  $X = (X_1, X_2, \dots)^\top$ . We define a bandwidth matrix  $H := \text{diag}(\underline{h}) = \text{diag}(h_1, h_2, \dots) \in \mathbb{R}^{\infty \times \infty}$ . We shall require that a norm  $\|\cdot\|$  can be admitted to the *weighted regressor* values and the *weighted point*, and for this the bandwidth sequence must be chosen appropriately. In particular, we shall assume that

$$H = hD = h \times \text{diag}(\phi_1, \phi_2, \dots), \quad (13)$$

where  $D \in \mathbb{R}^{\infty \times \infty}$  and  $h \in \mathbb{R}$ . By Kolmogorov's three-series theorem, the sequence of weighted regressors  $\{\phi_j^{-1}X_j\}$  is square summable, with probability one, provided that the marginal regressors  $X'_j$  are independent with finite variance and satisfy

$$\sum_{j=0}^{\infty} E \min \{1, \phi_j^{-2}X_j^2\} < \infty, \quad (14)$$

so that  $(\phi_1^{-1}X_1, \phi_2^{-1}X_2, \dots)^\top =: Z$  is  $(\ell_2, \|\cdot\|_2)$ -valued. In terms of the autoregressive framework the sequence  $\phi_j$  can be interpreted as non-decreasing weights that represent the “relative influence” of the marginal regressors, which diminishes as lags get further apart.

For this purpose we assume from now on that *the bandwidth-weighted*  $X$  and  $x$  (i.e.  $Z$  and  $z := (\phi_1^{-1}x_1, \phi_2^{-1}x_2, \dots)^\top$ , respectively) are  $\ell_2$ -valued<sup>3</sup> and normed with  $\|\cdot\| = \|\cdot\|_2$ . In view of this assumption, (with an abuse of notation) we can extend the usual definition of shifted small deviation to account for the generalized support  $[0, \lambda]$  and bandwidth vector  $\underline{h} = (h_1, h_2, \dots)^\top$ :

$$\varphi_x(\underline{h}\lambda) := P(\|H^{-1}(x - X_t)\| \leq \lambda) = P(\|D^{-1}(x - X_t)\| \leq h\lambda). \quad (15)$$

Equivalently,  $\varphi_x(\underline{h}\lambda) = P(X_t \in \mathcal{E}(x, \underline{h}\lambda))$ , where  $\mathcal{E}$  is the infinite-dimensional hyperellipsoid centred at  $x \in \mathbb{R}^\infty$ , and  $\lambda$  is as defined in section 2.2. Clearly,  $\varphi_x(\underline{h}\lambda) = \varphi_z(h\lambda)$ .

For later reference, we also define the joint small ball probability of the regressor vectors observed at different times  $t$  and  $s$  as the joint distribution

$$\psi_x(\underline{h}\lambda; t, s) := P((X_t, X_s) \in \mathcal{E}(x, \lambda\underline{h}) \times \mathcal{E}(x, \lambda\underline{h})). \quad (16)$$

### 3 The Estimator

We observe a sample  $\{Y_t, X_t\}_{t=1}^n$  with  $Y_t \in \mathbb{R}$  and  $X_t \in \mathbb{R}^\infty$ . We propose to estimate  $m(x) = E(Y|X = x)$ ,  $x \in \mathbb{R}^\infty$  by the following local constant type estimator:

$$\hat{m}(x) := \frac{\sum_{t=1}^n \mathcal{K}(H^{-1}(x - X_t)) Y_t}{\sum_{t=1}^n \mathcal{K}(H^{-1}(x - X_t))} \equiv \frac{\sum_{t=1}^n K(\|H^{-1}(x - X_t)\|) Y_t}{\sum_{t=1}^n K(\|H^{-1}(x - X_t)\|)}. \quad (17)$$

In the static case we may observe an infinity of regressors, but in the autoregression case we essentially observe only  $\{Y_1, Y_2, \dots, Y_n\}$  rather than the full infinity, see Assumptions A below. Hence for practical applications, one may employ a truncation argument on the regressor (as will be done in section 3.4 - albeit with a different purpose) and let the effective dimension  $\tau$  of the regressor  $X_t$  to increase in  $n$ .

The estimator can be viewed as an infinite-dimensional generalization of the standard multivariate local linear estimator, and is a special case of the one in Ferraty and Vieu (2002), Masry (2005) and references therein for functional data. In the following section we will examine some asymptotic properties of the estimator.

---

<sup>3</sup>This gives a mild restriction on the range of possible points at which the estimation is made; i.e.  $x \in \mathbb{R}^\infty$  is such that  $\sum_j j^{-2p} x_j^2 < \infty$ .

## 4 Asymptotic Properties

In this section we introduce the main results of our paper, deriving some large sample asymptotics of the proposed estimator (17). We establish consistency in both the pointwise and uniform sense, and also asymptotic normality. All proofs are detailed in the appendix.

We consider two different cases: (1) the static regression and (2) the dynamic regression. Below we specify two sets of dependence conditions, either of which will be assumed on the data generating process of the sample observations. Assumption A1 corresponds to the static regression case where we have exogenous regressors that are jointly observed in time in a weakly dependent manner. No restriction is needed as regards the dependence structure between the marginal regressors, although certain additional conditions can be potentially imposed at the later stage (see Assumptions C below). The second option A2 concerns the dynamic regression framework. In this case, the notion of near epoch dependence is adopted to describe the dependence structure of the processes defined as functions of the response variables. The assumptions below suggest that there is a trade-off between the degree of mixing and the possible order of moments,  $2+\delta$ , we allow on the response variable.

### ASSUMPTIONS A

- A1. *The marginal regressors  $X_{1t}, X_{2t}, X_{3t}, \dots$  are exogenous variables, and the sample data  $\{Y_t, X_t\}_{t=1}^n = \{Y_t, (X_{1t}, X_{2t}, \dots)\}_{t=1}^n$  is stationary and jointly arithmetically  $\alpha$ -mixing with rate  $k > (2\delta + 4)/\delta$ , where  $\delta$  is as defined in Assumption B4 below.*
- A2. *Each regressor is either a lag of the response variable  $Y_t$  or of a covariate  $V_t$ , i.e.  $X_{jt} = Y_{t-j}$  or  $X_{jt} = V_{t-j}$ ,  $j \in \mathbb{N}$ , and  $\{Y_t, V_t\}_{t=1}^n$  is stationary and arithmetically  $\alpha$ -mixing with rate  $k > (2\delta + 4)/\delta$ . Also, the process  $K_t := K(\|H^{-1}(x - X_t)\|)$  is near epoch dependent on  $(Y_t, V_t)$ , and there exists some  $r = r_n \rightarrow \infty$  such that the rate of stability for  $K_t$  denoted  $v_2(r_n) = v_2(r)$  satisfies*

$$v_2(r)^{1/2} [\varphi_x(\underline{h}\lambda)]^{-(2\delta+3)/(2\delta+2)} n^{1/(2(\delta+1))} \rightarrow 0 \quad \text{as } n \rightarrow \infty. \quad (18)$$

REMARK. Our model under Assumption A2 can be viewed as a generalization of the NAARX model in Chen and Tsay (1993). The framework nests both the fully autoregressive framework in which  $X_{jt} = Y_{t-j}$  for all  $j$ , and the case where the regressor vector consists only of the lags of a covariate  $V_t$ .

### 4.1 Pointwise consistency

Pointwise consistency of the local constant estimator was first studied by Watson (1964) and Nadaraya (1964) for i.i.d data with  $d = 1$ . Their result was extended to the multivariate case (finite dimension) by Greblicki and Krzyzak (1980) and Devroye (1981). Robinson (1983) and Bierens (1983) were amongst the earliest papers that worked on consistency of the estimator with dependent observations

(both static regression and autoregression were allowed in their frameworks), followed by Roussas (1989), Fan (1990), and Phillips and Park (1998) to name a few out of numerous papers. The case of the functional regressor was first studied by Ferraty and Vieu (2002).

In this section we establish the pointwise weak consistency of the estimator (17) with dependent data satisfying either A1 or A2. A set of assumptions required for the theory is now introduced, and some introductory arguments are briefly sketched.

## ASSUMPTIONS B

- B1. *The regression operator  $m : \mathbb{R}^\infty \rightarrow \mathbb{R}$  is continuous in some neighbourhood of  $x$*
- B2. *The marginal bandwidths satisfy  $h_j = h_{j,n} \rightarrow 0$  as  $n \rightarrow \infty$  for all  $j = 1, 2, \dots$ , where  $\text{diag}(h_1, h_2, \dots) = \text{diag}(\underline{h}) = H$  is the bandwidth matrix, and the small ball probability obeys  $n\varphi_x(\underline{h}\lambda) \rightarrow \infty$  for every point  $x \in \mathbb{R}^\infty$ , where  $\varphi_x(\underline{h}\lambda) := P(\|H^{-1}(x - X)\| \leq \lambda) \rightarrow 0$  as  $n \rightarrow \infty$ .*
- B3. *The kernel  $K$  is type-I*
- B4. *The response  $Y_t$  satisfies  $E(|Y_t|^{2+\delta}) \leq C < \infty$  for some  $C, \delta > 0$ .*
- B5. *The joint small ball probability (16) satisfies  $\psi_x(\underline{h}\lambda; i, j) \leq C\varphi_x(\lambda\underline{h})^2, \forall i \neq j$ .*
- B6. *The conditional expectation  $E(|Y_t Y_s| | X_t, X_s) \leq C < \infty$  for all  $t, s$ .*

REMARK. The continuity assumption B1 is necessary for asymptotic unbiasedness of the estimator. It will be shown that the estimator is unbiased at every point of continuity, and that the rate of convergence for the bias term can be specified upon imposing further smoothness condition on the regression operator, see later. Assumption B2 can be thought of as an extension of the usual bandwidth conditions that are assumed in finite-dimensional nonparametric literature, cf. (6). As discussed before,  $n\varphi_x(\underline{h}\lambda)$  can be understood as an approximate number of observations that are “close enough” to  $x$ . Therefore, it is sensible to postulate that  $n\varphi_x(\underline{h}\lambda) \rightarrow \infty$  as  $n \rightarrow \infty$ , meaning that the point  $x$  is visited many times by the sample of data as the size of the sample grows to infinity. This is in line with the usual assumption that  $nh^d \rightarrow \infty$  when  $X \in \mathbb{R}^d$ , in which case the small ball probability is given by  $\varphi_x(h) \propto h^d p_X(x)$  as noted in (6). Conditions B5 and B6 are imposed to control the asymptotics of the covariance terms. The validity of condition B5 can be easily seen in the  $\mathbb{R}^d$  frameworks; for relevant discussions, see Ferraty and Vieu (2006, Remark 11.2).

To sketch the idea, we write  $K_t := K(\|H^{-1}(x - X_t)\|)$  for the sake of simplicity of presentation (note its dependence upon  $X_t$ ), and express the estimator (17) as

$$\hat{m}(x) := \frac{\sum_{t=1}^n K(\|H^{-1}(x - X_t)\|) Y_t}{\sum_{t=1}^n K(\|H^{-1}(x - X_t)\|)} = \frac{\frac{1}{n} \sum_{t=1}^n \frac{K_t}{EK_1} Y_t}{\frac{1}{n} \sum_{i=1}^n \frac{K_t}{EK_1}} = \frac{\hat{m}_2(x)}{\hat{m}_1(x)}. \quad (19)$$

We then employ the following decomposition:

$$\begin{aligned}\widehat{m}(x) - m(x) &= \frac{\widehat{m}_2(x)}{\widehat{m}_1(x)} - m(x) = \frac{\widehat{m}_2(x) - m(x)\widehat{m}_1(x)}{\widehat{m}_1(x)} \\ &= \frac{E\widehat{m}_2(x) - m(x)E\widehat{m}_1(x)}{\widehat{m}_1(x)} + \frac{[\widehat{m}_2(x) - E\widehat{m}_2(x)] - m(x)[\widehat{m}_1(x) - E\widehat{m}_1(x)]}{\widehat{m}_1(x)},\end{aligned}\quad (20)$$

where clearly  $E\widehat{m}_1(x) = 1$ . Below we show consistency by proving that the ‘bias part’  $E\widehat{m}_2(x) - m(x)$  and the ‘variance part’  $[\widehat{m}_2(x) - E\widehat{m}_2(x)] - m(x)[\widehat{m}_1(x) - 1]$  are both negligible in large samples. As for the latter term, it suffices to show the mean squared convergence of  $\widehat{m}_2(x) - E\widehat{m}_2(x)$  to zero because  $\widehat{m}_1(x) \xrightarrow{P} 1$  then readily follows.

**THEOREM 1.** *Suppose that Assumptions B1-B5 hold. Then the estimator (17) with sample observations  $\{Y_t, X_t^\top\}_{t=1}^n$  satisfying either A1 or A2 is weakly consistent for the regression operator  $m(x)$ . That is, as  $n \rightarrow \infty$*

$$\widehat{m}(x) \xrightarrow{P} m(x). \quad (21)$$

In the following section, we present the rates of convergence and asymptotic normality under additional regularity conditions.

## 4.2 Asymptotic Normality

Earlier studies on the limiting distribution of the standard Nadaraya-Watson estimator can be traced back to Schuster (1972) and Bierens (1987), where the case of univariate and multivariate regressors was considered, respectively. The case of dependent samples was studied in Robinson (1983) and Bierens (1983), Masry and Fan (1997), and by many others under various model setups and different regularity conditions. Masry (2005, Theorem 4) and Delsol (2009) established general distribution theories for Nadaraya-Watson type estimators in a semi-metric space. Our results are different from those in two respects. First, the difference of our framework from the functional literature discussed in the beginning of Section 2 gives us some additional flexibility. Second, whereas the final results of many existing papers were given in terms of abstract functions, our results are presented with an explicit rate of convergence.

The primary objective of this section is to outline the main theory and some interesting consequences thereof. Both cases of the independent marginals (in other words, when the marginal regressors  $X_j$  are independent and identically distributed) and also a dependent framework are allowed. Specifically, we introduce how independence restriction can possibly be moderated to allow for some mild dependence structure. In particular, the second condition in Assumption C below specifies the extent to which certain cross-sectional dependence can be allowed on the marginal regressors in our theory while still allowing for specification of the exact form of the convergence rate of the estimator.

ASSUMPTIONS C. *For every fixed  $t$ , the real-valued stochastic process formed by the marginal regressors  $\{X_{jt}\}_{j=1}^{\infty}$  is either:*

C1. *independent and identically distributed over  $j$  with  $EX_{jt}^4 \leq C < \infty \forall j$ , or*

C2. *stationary (over  $j$ ), and admits a moving average representation:*

$$X_{jt} = \sum_{u=-\infty}^{\infty} a_u \epsilon_{j-t-u}, \quad (22)$$

*where  $a_u$  is a square summable sequence, and  $\{\epsilon_{jt}\}_j$  is an independent and identically distributed standard Gaussian sequence.*

REMARK. In either case the marginal regressors are required to be identically distributed over  $j$ ; an additional distributional assumption will be imposed in D2 below. Nonetheless, the possible degree of dependence allowed in C2 is very mild and general, since an equivalent condition of having the representation for a Gaussian process is simply the existence of the spectral density. Note that (22) includes the causal (one-sided) MA representation as a special case. If a stationary stochastic process  $\{X_{jt}\}_j$  is  $\alpha$ -mixing (over  $j$ ), then it always has such a representation (i.e.  $a_u = 0, \forall u < 0$ ) provided it is Gaussian. This is because any  $\alpha$ -mixing process is regular<sup>4</sup> by definition, so is linearly regular when it is Gaussian, and hence (with stationarity) admits the Wold decomposition with independent Gaussian innovations by Corollary 17.3.1 of Ibragimov and Linnik (1971).

Note that each C1 and C2 is consistent with the case allowed in Assumption A1 and A2, respectively (because in the latter case the process  $\{X_{jt}\}_j$  consists of temporal lags of the response variable and/or a covariate which form a mixing process by Assumption A2), although the dependence structure specified in C2 can be allowed also for the static case (i.e. A1). This suggests that there is absolutely no need to assume independence between marginal regressors in our model (2) under Gaussianity, and hence a wide flexibility is allowed in terms of the model setup. In particular, the convergence rates of our estimator will be shown to be the same (upto some constant factor in the limiting variance) in both cases C1 and C2. Lastly, the requirement of a finite 4th moment is imposed to ensure that the squared marginal regressors have finite second moments due to the reasons to be clarified below; obviously, when a lag of the response is included in the dynamic regression framework (A2), this forces  $\delta \geq 2$  in Assumption B4.

We now introduce some main assumptions needed for distributional theories.

---

<sup>4</sup>In the sense of Ibragimov and Linnik (1971) and Davidson (1994, Part III)



### 4.2.1 The ‘bias component’

The first part concerns with the asymptotic ‘bias’, where Assumptions A is strengthened by imposing additional smoothness conditions and suitable bandwidth adjustments. They belong to a set of sufficient conditions under which the exact upper bound of the asymptotic bias can be specified. Note that alternatively, a Fréchet-type differentiability condition can be imposed, as was done in Mas (2012).

#### FURTHER ASSUMPTIONS B

B7. *The regression operator  $m : \mathbb{R}^\infty \rightarrow \mathbb{R}$  satisfies*

$$|m(x) - m(x')| \leq \sum_{j=1}^{\infty} c_j |x_j - x'_j|^\beta \quad (23)$$

*for every  $x, x' \in \mathbb{R}^\infty$ , and some constant  $\beta \in (0, 1]$ , where  $\{c_j\}$  is some sequence of real constants that satisfies  $\sum_{j=1}^{\infty} c_j \leq 1$ .*

B8. *The marginal bandwidths satisfy  $h_j = \phi_j \cdot h$  for some positive real numbers  $\phi_j$ , where  $h = h_n \rightarrow 0$  as  $n \rightarrow \infty$ . We suppose that  $\phi_j$  satisfy  $\sum_{j=1}^{\infty} \phi_j^{-2} < \infty$  and  $\sum_{j=1}^{\infty} c_j \phi_j^\beta < \infty$ .*

REMARK. These additional assumptions help us specify and regulate the bias component. Assumption B8 extends the previous bandwidth condition B2. Obviously, it is consistent with what was previously assumed in B2 since  $h \rightarrow 0$  implies the coordinate-wise convergence of each marginal bandwidths. With this condition one is able to write the asymptotic bias expression and the order of the bias-variance balancing bandwidth in terms of the common factor  $h$ . It is possible to dispense with this condition at the cost of imposing minor modifications in B7; the asymptotic bias will then be written in terms of the infinite sum of a weighted marginal bandwidth  $h_j$ , whose convergence needs to be ensured. For the sake of understanding the asymptotic behaviour of the variance component, a further increment condition will be imposed on the sequence of marginal coefficients  $\phi_j$  in Assumption D later. We remark that at this point such an assumption is not necessary as the variance term is not concerned.

Assumption B7 replaces and strengthens Assumption B1, and can be thought of as a variant of Hölder-type continuity; the case of  $c_j = 2^{-j}$  and  $\beta = 1$  is implied by the Lipschitz condition. Another example of  $c_j$  includes  $\exp(-j)$ . Indeed, under B7 the regression operator becomes a contraction mapping, and the contribution from each marginal dimension decreases in lag or index. This ensures summability of the bias of the estimator and allows one to specify its order of convergence rate, cf. (29) below.

In the context of autoregression where  $X_j \equiv Y_{t-j}$  for all  $j$ , the model is given by

$$Y_t = m(Y_{t-1}, Y_{t-2}, \dots) + \varepsilon_t \quad (24)$$

and whether the stationary solution  $\{Y_t\}$  indeed exists is an important question, since (24) essentially gives an infinite number of recurrence relations whose solution may not be always well-defined. In the study of a class of general nonlinear AR( $d$ ) models, Duflo (1997) and Götze and Hipp (1994) assumed what is called the Lipschitz mixing condition (or the strong contraction condition), which is essentially (23) replaced by finite  $d$ -sum on the right hand side. In our context, Assumption B7 plays a similar role; Doukhan and Wintenberger (2008) showed that (23) with  $\sum_{j=1}^{\infty} c_j < 1$ , is sufficient for the existence of a stationary solution: for some measurable  $f$ ,

$$Y_t = f(\varepsilon_t, \varepsilon_{t-1}, \dots), \quad (25)$$

where  $\varepsilon_t$  is an i.i.d. sequence. Wu (2011) arrived at the same conclusion under the assumption of  $\sum_{j=1}^{\infty} c_j = 1$ ; the specific restrictions on  $c_j$  are chosen to reflect their findings, despite the fact that we are not restricting the error process  $\{\varepsilon_t\}$  to be an independent sequence in our model setup.

Before we proceed, we remark that from now on the rate condition stipulated in (18) is slightly strengthened as follows (and Assumption A2 is modified accordingly):

$$v_2(r)^{1/2}[\varphi_x(\underline{h}\lambda)]^{-1}n^{1/2} \rightarrow 0 \quad \text{as } n \rightarrow \infty. \quad (26)$$

#### 4.2.2 The ‘variance component’

We now move on to the second chunk of assumptions that are concerned with the ‘variance part’. As before, vectors  $Z$  and  $z$  are taken to mean  $(\phi_1^{-1}X_1, \phi_2^{-1}X_2, \dots)^\top$  and  $(\phi_1^{-1}x_1, \phi_2^{-1}x_2, \dots)^\top$ , respectively, where the vector  $x = (x_1, x_2, \dots)^\top$  is the point at which estimation is made, and  $\phi'_j$ s are the coefficients in Assumption B8.

#### ASSUMPTIONS D

- D1. *The induced probability measure  $P_{z-Z}$  is dominated by the measure  $P_Z$ , and its Radon-Nikodym density  $dP_{z-Z}/dP_Z =: p^*$  is continuous and is bounded away from zero at  $0 \in \mathbb{R}^\infty$ ; i.e.,  $p^*(0) > 0$ .*
- D2. *The distribution  $F$  of  $X_s^2$ , where each  $X_s$  is the marginal regressor, is regularly varying near zero with strictly positive index  $(-\rho) > 0$ .*
- D3. *Further to B8, the bandwidth satisfies  $h_j = j^p h$  (i.e.  $\phi_j = j^p$ ) with  $p \in \Pi(c, \beta)$ , where*

$$\Pi(c, \beta) = \left\{ p : \sum_{j=1}^{\infty} c_j j^{p\beta} < \infty, p > 1/2 \right\}.$$

- D4. *The conditional variance  $\text{var}[Y_t|X_t = u] = \sigma^2(u)$  is continuous in some neighbourhood of  $x$ ; i.e.  $\sup_{u \in \mathcal{E}(x, \underline{h}\lambda)} [\sigma^2(u) - \sigma^2(x)] = o(1)$ . Similarly, the cross-conditional moment  $E[(Y_t - m(x))(Y_s - m(x))|X_t = u, X_s = v] = \sigma(u, v)$ ,  $t \neq s$  is continuous in some neighbourhood of  $(x, x)$ .*

D5.  $R_{nt} := (EK_1)^{-1}\{K_t(Y_t - m(x)) - EK_t(Y_t - m(x))\}$  belongs to the domain of attraction of a normal distribution.

REMARK. Assumption D1 is concerned with a transition of the shifted small ball probability to the centred small deviation (whose asymptotic behaviour is more accessible), see Mas (2012). The explicit form of the derivative (and hence of the relationship between the two probabilities) cannot be easily computed in general. Nonetheless, in the special case of the Gaussian process  $Z$  with covariance operator  $\Sigma$  it is known by Sytaya (1974) and Zolotarev (1986) that

$$P(\|z - Z\| \leq \epsilon) \simeq P(\|Z\| \leq \epsilon) \exp \left\{ -\frac{1}{2} \|\Sigma^{-1/2} z\|^2 \right\} \quad \text{as } \epsilon \rightarrow 0. \quad (27)$$

The reader is directed to Li and Shao (2001) for detailed discussion on this asymptotic equivalence relation. Note that  $\Sigma$  can be expressed in terms of the  $a_j$  constants (in Assumption C), which govern the dependence between the marginal regressors and the bandwidth weights  $\phi_j$ :

$$\text{cov}(Z) = \Sigma = (DA)(A^*D), \quad (28)$$

where  $A = (a_{ij}) = (a_{i-j})$  and  $D = \text{diag}(\phi_1, \phi_2, \dots)$ .

Condition D2 is equivalent to saying that

$$\lim_{x \rightarrow \infty} \frac{F(1/(\gamma x))}{F(1/x)} = \gamma^\rho,$$

where  $\rho$  is the index of variation which is strictly negative. Under this condition, Dunker, Lifshits and Linde (1998, cf. Conditions *I* and *L*) derived the explicit behaviour of the small ball probability. We require the function  $F(1/x)$  to be regularly varying in order to ensure that the small ball probability is *well-behaved* near infinity in the asymptotic sense. Since only those functions having strictly negative  $\rho$  satisfy the condition, the distribution  $F$  of the squared regressor must be such that  $F(1/x)$  decreases (as  $x \rightarrow \infty$ ) at a *reasonable speed*. By reasonable we mean that the relative weight of decrease follows a power law, and the variation should be continuous. A large class of common distributions satisfies this condition; for example: the Gamma, Beta, Pareto, Exponential, Weibull, and also the Chi-squared distribution (in which case each  $X_s$  is Gaussian). Indeed, both D1 and D2 hold under Gaussianity (e.g. when condition C2 is assumed).

The specific bandwidth increment condition assumed in D3 is one framework under which the explicit behaviour of the small ball probability can be specified (cf. Dunker et al. (1998)). In the exceptional case of static regression where the regressors form an i.i.d. sequence, the probability can also be derived when the weights are of an exponential type (i.e.  $h_j = e^j h$ ) up to an unknown function, or are non-increasing in a particular manner (cf. Gao et al. (2003)) similar to the polynomial decay. In this paper however, we shall confine our attention to the case of the polynomial law for expositional simplicity and consistency of presentation, since the asymptotic behaviour of the small

ball is not yet known in the dependent case for choices other than the polynomial decay as in D3. In practice, we would require some ordering for the marginal regressors in the static regressions case A1, since the influence of marginals is set to decrease via the bandwidth adjustments. The standard conditions in D4 are assumed to deal with the asymptotics of the variance and covariance terms. The last condition is imposed to establish the self-normalized CLT without assuming higher moment conditions; relevant discussions can be found for example in de la Peña et al. (2009). The condition is not affected by the dependence structure of  $R_{nt}$  as the property is inherited to the approximated sum in the Bernstein's blocking procedure; see (77) for details.

With reference to (20) we are now able to derive the following results for the bias and variance components using Assumptions B7, B8 and C, and Corollary 1:

$$\mathcal{B}_n(x) := \left[ E\widehat{m}_2(x) - m(x) \right] \leq h^\beta \lambda^\beta \sum_{j=1}^{\infty} c_j j^{p\beta} \quad (29)$$

$$\mathcal{V}_n(x) := \text{var}[\widehat{m}_2(x)] \simeq \frac{\sigma^2(x)\xi_2}{n\varphi_x(h\lambda)\xi_1^2}, \quad (30)$$

where  $\lambda$  and  $\widehat{m}_2(\cdot)$  are as in (4) and (19), respectively. Formal derivation is done in 7.2 of the appendix. We next present the CLT of our estimator.

#### 4.2.3 Limiting distribution under independence of regressors

We first consider the situation in which there is a set of independent exogenous regressors in the static regression context. That is, when marginal regressors  $X_s$  are independent to each other and are identically distributed (i.e. satisfies Assumption C1), and the sample observations follow Assumption A1.

In this case, the asymptotic normality can be established for regressors that follow a wide range of different distributions. Recall that under Assumption D2, the distribution function  $F$  (of  $X^2$ ) is regularly varying with the index of variation  $\rho < 0$ . Then, by the characterization theorem of Karamata (1933) (see for example Feller (1971)), there always exists a slowly varying function  $\ell(x)$  that satisfies

$$F(1/x) = x^\rho \ell(x). \quad (31)$$

Now fix some  $p$ , the order of increment constant for bandwidth in Assumption D3, and denote by  $\mathcal{L}(t)$  the Laplace transform of  $X^2$ . We then define the following constants:

$$C_\ell = \lim_{\delta \rightarrow 0} \left[ \ell^{-1/2} \left( \delta^{-\frac{4p}{2p-1}} \right) \right], \quad \zeta = - \int_0^\infty \frac{u^{-1/2p} \mathcal{L}'(u)}{\mathcal{L}(u)} du$$

$$C^* = \frac{(2\pi)^{(1+2p\rho)}(2p-1)}{\Gamma^{-1}(1-\rho) \cdot (2p)^{\frac{2p(\rho+2)-1}{2p-1}}} \cdot \zeta^{\frac{2p(1+\rho)}{2p-1}}, \quad C^{**} = (2p-1) \cdot \left( \frac{\zeta}{2p} \right)^{2p/(2p-1)}$$

$X_j^2 \sim F$ i.i.d.	$\rho$	$\lim_{x \rightarrow \infty} \ell(x) = C_\ell^{-2}$	$\zeta$
Uniform(1,b)	-1	1	n/a
Gamma( $\alpha, \beta$ )	$-\alpha$	$\beta^\alpha \alpha^{-1} \Gamma(\alpha)^{-1}$	$\frac{\alpha \pi \beta^{-1/2p}}{\sin(\pi/2p)}$
exp( $\eta$ )	-1	$\eta$	$\frac{\pi \eta^{-1/2p}}{\sin(\pi/2p)}$
Weibull( $\alpha, \beta$ )	$-\alpha$	$\beta$	n/a
Pareto( $\theta, \mu$ )	-1	$\mu/\theta$	n/a
$\chi_1^2$	-1/2	$(2/\pi)^{1/2}$	$\frac{\pi 2^{(1-2p)/2p}}{\sin(\pi/2p)}$

Table 1: Examples of key constants for some common distributions

$$\kappa_0(K, p, F) = C^{**} \lambda^{-\frac{2}{2p-1}} \quad \text{and} \quad \kappa_1(K, p, F) = \frac{C^* C_\ell \xi_2}{p^*(0) \xi_1^2 \lambda^{\frac{1+2\rho p}{2p-1}}},$$

where  $\Gamma(\cdot)$  is the Gamma function,  $\xi_1$  and  $\xi_2$  are the constants specified in (11) (which simplify in case of uniform kernel for example),  $\lambda$  is the upper bound of the support of the kernel, and  $p^*(\cdot)$  is the Radon-Nikodym derivative in D1. The underlying arguments for the formulation of these constants can be found in Dunker, Lifshits and Linde (1998). To aid the exposition, we compute the constants for some common, regularly varying distributions in Table 1. The main result of this subsection now follows. The theorem gives the limiting distribution of the infinite-dimensional Nadaraya-Watson type estimator under cross sectional independence with respect to mixing sample data.

**THEOREM 2.** *Suppose that B2-B8 and D1-D4 hold. Let the marginal regressors  $X_s$  satisfy Assumption C1. Then the estimator (17) based on the sample observations  $\{Y_t, X_t\}_{t=1}^n$  satisfying A1 is asymptotically normal with the following limiting distribution:*

$$\sqrt{nh^{\frac{1+2\rho p}{2p-1}} \exp\left(-\kappa_0 h^{-\frac{2}{2p-1}}\right)} \left(\widehat{m}(x) - m(x) - \mathcal{B}_n(x)\right) \Rightarrow N\left(0, \kappa_1 \sigma^2(x)\right), \quad (32)$$

where  $\mathcal{B}_n(x) = O(h^\beta)$  is the bias component as in (29) and  $\sigma^2(\cdot)$  is the conditional variance defined in Assumption D4.

#### 4.2.4 Limiting distribution under Gaussianity & dependence of regressors

The independence condition between the regressors assumed in the previous section can be relaxed to allow some mild dependence specified in Assumption C2. In doing so, we make use of the result derived in Hong, Lifshits and Nazarov (2016, Theorem 1.1), where the asymptotics of the small deviation probability of Gaussian dependent sequences was investigated. This setting not only grants sufficient flexibility in the static regression case, but moreover allows one to compute the distributional result for the dynamic regression context, where the regressor vector consists of time lags of the

response or a covariate with dependence structure stipulated in Assumption A2. The price we have to pay for this modification is the Gaussianity restriction on the regressors.

With reference to Table 1 above, we can easily compute the constants  $C^*$  and  $C^{**}$  for the Gaussian case, denoted  $C_G^*$  and  $C_G^{**}$  respectively, as follows:

$$C_G^* = \frac{(2\pi)^{(1-p)}(2p-1)}{2 \cdot (2p)^{\frac{3p-1}{2p-1}}} \cdot \left[ \frac{\pi 2^{(1-2p)/2p}}{\sin(\pi/2p)} \right]^{\frac{-p}{2p-1}}, \quad C_G^{**} = \frac{2p-1}{2} \left( \frac{\pi}{2p \sin \frac{\pi}{2p}} \right)^{\frac{2p}{2p-1}}.$$

For the square summable sequence  $a_j$  in (22) define

$$C_{\mathcal{A}} = \left[ \frac{1}{2\pi} \int_0^{2\pi} \left| \sum_{j=0}^{\infty} a_j \exp(ijs) \right|^{1/p} ds \right]^p \quad \text{and} \quad \kappa_2(K, p, a) = \frac{C_G^* C_{\ell} \xi_2 \xi_1^{-2}}{e^{-\frac{1}{2} \|\Sigma^{-1/2} z\|_2^2} (C_{\mathcal{A}} \lambda)^{\frac{1-p}{2p-1}}}.$$

where  $z = (z_j) = (j^{-p} x_j) = D^{-1}x$ . Recall that for the uniform (Box) kernel  $\xi_2 = \xi_1^2$ , so they cancel out in  $\kappa_2$ . Let  $\kappa'_0 = C_G^{**} \lambda^{-2/(2p-1)}$ .

With other constants defined as before, we now have the following asymptotic normality for the case of dependent regressors. We reiterate that the result covers both the static and dynamic regressions context (A1 and A2), and is invariant to (32) upto a constant factor in the limiting variance as long as the cross-dependence structure satisfies Assumption C2.

**THEOREM 3.** *Suppose B2-B8 and D1-D4 hold. Let the regressor  $X = (X_1, X_2, \dots)^\top$  is jointly normally distributed with zero mean and the covariance operator  $\Sigma$ , and satisfies C2. Then, the estimator (17) with respect to sample observations  $\{Y_t, X_t^\top\}_{t=1}^n$  satisfying either A1 or A2 is asymptotically normal with the following limiting distribution:*

$$\sqrt{nh^{\frac{1-p}{2p-1}} \exp\left(-\kappa'_0 h^{-\frac{2}{2p-1}}\right)} (\widehat{m}(x) - m(x) - \mathcal{B}_n(x)) \implies N(0, \kappa_2 \sigma^2(x)), \quad (33)$$

where  $\mathcal{B}_n(x)$  is the bias component in (29) and  $\sigma^2(x)$  is the conditional variance defined in Assumption D4.

**REMARK.** The additional constant  $C_{\mathcal{A}}$  is a function of the sequence  $a_j$ , and represents the dependence structure allowed between the regressors. This suggests an interesting finding that says allowing for dependence does not incur much penalty; we conjecture that similar conclusion would hold for regressors of different distributions, but leave it for future studies. The exponential term in the denominator of the asymptotic variance arises from the asymptotic equivalence relationship between the shifted and non-shifted small deviation for  $\ell_2$ -valued Gaussian variables, cf. (27).

In both frameworks of independent regressors and dependent Gaussian regressors we are able to construct the self-normalised central limit theorem; define

$$\Delta_n^2(x) = \sum_{t=1}^n \left( \sum_{s=1}^n K_s \right)^{-2} \left[ K_t (Y_t - \widehat{m}(x)) \right]^2, \quad (34)$$

where  $K_t := K(\|H^{-1}(x - X_t)\|)$  as before.

**COROLLARY 2.** *Further to the conditions assumed either in Theorem 2 or Theorem 3, suppose that Assumption D5 holds. Then the following central limit theorem holds*

$$\Delta_n^{-1}(x) \left( \widehat{m}(x) - m(x) - \mathcal{B}_n(x) \right) \implies N(0, 1),$$

where  $\Delta_n(x)$  is the square root of (34).

This self-normalized limit distribution gives (pointwise) confidence intervals for  $\widehat{m}(x)$ , which can be used as a basis for conducting standard statistical inference.

### 4.3 Optimal Bandwidth

We now briefly discuss the issue of bandwidth optimality. As in the finite-dimensional framework, there is a bias-variance trade-off. As the bandwidth goes up, the variance gets smaller while the bias increases, and vice versa. Therefore we search for the optimal bandwidth  $h_{opt}$  that balances the order of those two quantities.

We first suppose that  $p \in \Pi(c, \beta)$  is given. In the i.i.d. case with Gaussian regressor we have

$$h^\beta \sim \sqrt{\frac{\exp(\kappa_0 h^{-2/(2p-1)})}{n h^{\frac{1-p}{2p-1}}}}, \quad (35)$$

so that

$$\left[ 2\beta + \frac{1-p}{2p-1} \right] \cdot \log h - \kappa_0 h^{-\frac{2}{2p-1}} \sim -\log n.$$

Taking  $h \sim (\log n)^a$  for some  $a < 0$  balances the leading terms on both sides:

$$\left[ 2\beta + \frac{1-p}{2p-1} \right] \cdot a \cdot \log \log n - \kappa_0 (\log n)^{-\frac{2}{2p-1} \cdot a} \sim -\log n. \quad (36)$$

The explicit order  $a$  that solves (36) can be expressed in terms of  $n$ ,  $\beta$  and  $p$ . Writing  $\vartheta := [2\beta + (1-p)/(2p-1)]$  and  $\chi := 2/(2p-1)$  for notational simplicity, and solving for  $a$  we have

$$a_{opt} = \frac{\vartheta \cdot \mathcal{W}\left(\frac{\chi}{\vartheta} \cdot \kappa_0 \cdot n^{\chi/\vartheta}\right) - \chi \log n}{\vartheta \chi \cdot \log \log n}, \quad (37)$$

where  $\mathcal{W}(y)$  is the Lambert W function (see e.g. Olver et al. (2010)), which returns the solution  $x$  of  $y = x \cdot e^x$ . From (37) the optimal bandwidth  $h_{opt} \sim (\log n)^{a_{opt}}$  follows in which case the asymptotic root mean squared error is of the order  $(\log n)^{\beta a_{opt}}$ .

**REMARK.** We can look for the optimal bandwidth for the cases of non-Gaussian regressors by following exactly the same manner as above; tedious details are omitted here. As regards the solution in (37), since the mapping  $x \mapsto x \cdot e^x$  is not an injection, the solution may be multi-valued on

the negative domain, i.e.  $y < 0$ . This does not happen in (37) provided  $\beta \geq 1/4$  (however big  $p$  is), because  $(1-p)/(2p-1)$  is bounded away from  $-1/2$ ; in this case, the coefficient of the double logarithmic term in (36) is strictly smaller or equal to zero.

Since the log terms dominate the double logarithm in (36) as the sample size  $n$  increases, it can be readily expected that the optimal value of  $a$  in (37) converges to a limit in such a way that the leading orders are balanced. Below we introduce without formal justification a trivial result that gives the lower bound (infimum) of the optimal bandwidth (and hence of the optimal rate that balances the bias and variance, see also Mas (2012, Theorem 3)). It is worth noting that the result below holds for other choices of the distribution of the regressors, since the exponent of the leading term  $-2/(2p-1)$  remains invariant as it was shown in (32) and (33).

**COROLLARY 3.** *For any fixed choice of  $p \in \Pi(c, \beta)$  and the distribution  $F$  of  $X^2$  satisfying D2, the order of the optimal bandwidth  $a_{opt}$  satisfies*

$$a_{opt} \downarrow \left( -\frac{2p-1}{2} \right) \quad \text{as } n \rightarrow \infty, \quad (38)$$

*which suggests that the lower bound of the optimal bandwidth is given by*

$$(\log n)^{-\frac{2p-1}{2}} \preceq h_{opt} \sim (\log n)^{a_{opt}}. \quad (39)$$

**REMARK.** (i) This result tells us what the best possible performance we can expect from the optimal bandwidth. Because  $n^k(\log n)^{-(2p-1)/2} \rightarrow \infty$  for any positive real number  $k$ , it follows that we cannot possibly estimate the regression function at a polynomial rate. (ii) The above arguments are true for any  $p \in \Pi(c, \beta)$ . Let  $p_{\max} = \sup_{p \in \Pi(c, \beta)} p$ . Then a lower bound on the optimal rate of convergence (over all  $p$ ) is  $(\log n)^{-(p_{\max}-\frac{1}{2})}$ . For example, when  $c_j = (1/2)j^{-2}$  we have  $p_{\max} = 1/\beta$ . Unfortunately, it is generally the case that  $p_{\max} \notin \Pi(c, \beta)$ , in which case the lower bound is not quite achievable by our method.

**REMARK.** Regarding bandwidth selection, one possibility is the Bayesian bandwidth selection methods like proposed in Zhang, King, and Hyndman (2006). We take as prior for  $h$  the density proportional to  $1/(1+\lambda h^2)$  and as prior for  $p-1/2$  the density of a  $\chi^2(w)$  random variable. The hyperparameters  $\lambda, w$  may be chosen by experimentation. The priors are combined with a Gaussian (least squares) density to deliver a posterior for the bandwidth.

We conclude this section by briefly considering the case where we have an exponential growth in the bandwidth. That is, when  $\phi_j = \exp(j^q)$ , so that  $h_j = \exp(j^q)h$  for some  $q > 0$ , and let  $q \in \Phi(c, \beta) = \{q : \sum_{j=1}^{\infty} c_j(\exp(j^q))^{\beta} < \infty\}$ . We will only briefly go over this because explicit



expressions for the small ball probability are limited in this case. Specifically, it is known that the small ball probability is of the following order

$$\begin{aligned}\varphi_x(\underline{h}) &:= P(\|H^{-1}(x - X_t)\| \leq 1) \\ &= P(\|D^{-1}(x - X_t)\| \leq h) \sim \exp\left[-(\log(1/h))^{1+1/2q}\right],\end{aligned}\tag{40}$$

provided that the regressor is Gaussian and they are cross-sectionally independent (Assumption C1). No such explicit expression is available when the regressor follows other distributions, Dunker et al. (1988), and also, it is not known whether this still holds under cross-sectional dependence (e.g. under Assumption C2). Now it follows that the optimal bandwidth  $h_{opt}$  that balances the squared bias and the variance satisfies the following equivalence relation:

$$h^{2\beta} \sim \frac{1}{n \exp\left[-(\log(1/h))^{1+1/2q}\right]}$$

from which we have

$$\log n + 2\beta \log h \sim \left[\log\left(\frac{1}{h}\right)\right]^{1+1/2q}.$$

Try  $h = \exp[a \cdot (\log n)^b]$  for some  $a < 0$  and  $b > 0$ , and solve for  $n$ , then we have

$$[\log n + 2\beta a (\log n)^b]^{\frac{2q+1}{2q}} \sim -a (\log n)^b$$

and the optimal bandwidth is the one with the values  $a$  and  $b$  that satisfies this relation. Now in view that  $\beta$  is non-negative and  $a < 0$ , it can be shown that

$$\exp\left[-(\log n)^{\frac{2q}{2q+1}}\right] \leq h_{opt} \sim \exp[a_{opt} \cdot (\log n)^{b_{opt}}]$$

The performance is better in general than what we had with a polynomial increment for bandwidth, although we still do not attain polynomial rate of convergence.

## 4.4 Uniform consistency

Uniform consistency of the Nadaraya-Watson estimator was first studied by Nadaraya (1964, 1970) and subsequently by numerous others. To introduce a few, Devroye (1978) weakened the regularity conditions required in the previous papers, and Robinson (1983) proved uniform consistency for dependent sample data. In the functional statistics literature, only uniform consistency with respect to i.i.d. data sample has been established so far, see Ferraty et al. (2010). We start by introducing the notion of Kolmogorov's entropy.

**DEFINITION 5.** *Given some  $\eta > 0$ , let  $L(S, \eta)$  be the smallest number of open balls in  $E$  of radius  $\eta$  needed to cover the set  $S \subset E$ . Then Kolmogorov's  $\eta$ -entropy is defined as  $\log L(S, \eta)$ .*

From the definition it can be readily expected that Kolmogorov's entropy is dependent heavily on the nature of the space we work on, and is closely related to the rate of convergence of the estimator.

It is well known that the regression function cannot be estimated uniformly over the entire space, see e.g. Bosq (1996). In our infinite dimensional framework, even greater restrictions apply; since we are working on infinite sequence spaces, none of their subsets can be covered by a finite number of balls, so that  $L(S, \eta) = \infty$ . Therefore, we propose to adopt a truncation argument and consider uniform consistency over a set whose effective dimension is increasing in sample size  $n$ . In particular, we define the set

$$S_\tau := \{u | (u_i)_{i \in \mathbb{Z}^+}, u_j = 0 \text{ for all } j > \tau, \|u\|_\infty \leq \lambda\} \quad (41)$$

where  $\tau = \tau_n$  is some increasing sequence and  $\lambda$  is fixed, and consider uniform consistency over this compact set.

Then Kolmogorov's entropy of the set  $S_\tau$  is given as follows:

LEMMA 2. *Kolmogorov's  $\eta$ -entropy of  $S_\tau$  defined in (41) with  $\tau = \tau_n (\rightarrow \infty)$  and  $\lambda > 0$  is*

$$\log L(S_\tau, \eta) = \log \left[ \left( \frac{2\lambda\sqrt{\tau}}{\eta} + 1 \right)^\tau \right]. \quad (42)$$

REMARK. We note that (42) is indeed in line with common intuition; as the dimension  $\tau$  increases, the number of balls with some fixed radius required to cover the set goes off to infinity. The proof of this result can be done by exploiting the splitting technique and then by attempting to cover the polyhedron of increasing dimension. See appendix for details. From this result it follows that for fixed  $\lambda$  and  $\eta = \eta_n$ , Kolmogorov's entropy  $\log L(S_\tau, \eta)$  is of order  $O(\tau \log \tau - \tau \log \eta)$ .

We now introduce some further assumptions needed for uniform consistency:

#### ASSUMPTIONS E

E1. *For sufficiently large  $n$ , Kolmogorov's  $\eta$ -entropy  $\log L(S_\tau, \eta)$  satisfies*

$$\frac{(\log n)^{8+2\epsilon}}{n\varphi_x(\underline{h})} \leq \log L(S_\tau, \eta) \leq \frac{\sqrt{n\varphi_x(\underline{h})}}{(\log n)^{1+\epsilon}} \quad \text{for some } \epsilon \in (0, 1/2). \quad (43)$$

*Furthermore,  $0 < \varphi_x(\underline{h}) \preceq h < \infty$  and  $(\log n)^2 / (n\varphi_x(\underline{h})) \rightarrow 0$  as  $n \rightarrow \infty$ .*

E2. *The kernel function  $K$  is Lipschitz continuous on  $[0, \lambda]$ .*

REMARK. The first part of Assumption E1 specifies the rate at which Kolmogorov's entropy should behave with sample size  $n$  (hence in dimension  $\tau = \tau_n$ ). From the upper and lower bound it readily follows that  $n\varphi(h)$  must be of order larger than  $(\log n)^{6+2\epsilon}$ . This assumption allows sufficient generality; for example, the restriction that the bias-variance optimal bandwidth satisfies

$h \succeq (\log n)^{-(2p-1)/2}$  gives  $n\varphi(h) \preceq (\log n)^{(2p-1)\beta}$ . In this case, assumption (43) is valid as long as  $p$  is moderately large enough relative to  $\beta \leq 1$  in such a way that  $6 + 2\epsilon \leq (2p-1)\beta$ . The second part of E1 is standard and the last condition straightforwardly follows by (43) and only slightly strengthens the bandwidth condition in Assumption B2.

We now introduce the main result of this section. Note that in the sequel, (with a slight abuse of notation)  $X$  is taken to denote the regressor, but with zeros after its  $\tau^{th}$  ( $= \tau_n \rightarrow \infty$  as  $n \rightarrow \infty$ ) entry; that is,  $X = (X_1, X_2, \dots, X_\tau, 0, 0, \dots)^\top$  (so that the original  $X$  is recovered as  $n \rightarrow \infty$ ). Also, the regression operator and the estimator with respect to this truncated regressor are denoted by  $m_\tau(\cdot)$  and  $\hat{m}_\tau(\cdot)$ , respectively. The aforementioned assumptions are understood to be modified accordingly.

For uniform consistency we impose a slightly stricter condition on the response:

B4'. *The response  $Y_t$  satisfies the following tail condition: There exists some positive constant  $\gamma_1$  and  $C$  such that  $P(|Y_t| > u) \leq C \exp(1 - u^{\gamma_1})$  for any  $u > 0$ .*

For example, a Gaussian random variable satisfies B4' with  $\gamma_1 < 2$ . The condition is also satisfied by many unbounded variables and all those bounded ones as well. We also impose a stronger condition on mixing coefficients; from hereafter, by A1' and A2' we mean Assumptions A1 and A2 but with the arithmetic mixing rate condition strengthened to the following exponential mixing condition (cf. Definition 1):

$$\alpha(r) \leq \exp(-\varsigma r^{\gamma_2}) \quad (44)$$

where  $\varsigma > 1$  and  $\gamma_2$  is a positive constant such that  $\gamma := 1/(\gamma_1^{-1} + \gamma_2^{-1}) \geq 1$ . In the case of bounded response (i.e.  $|Y_t| \leq C$ ),  $\gamma_1$  may be taken to be  $\infty$  so that  $\gamma_2 = \gamma \geq 1$ .

**THEOREM 4.** *Suppose that Assumptions B2, B3, B4', B5-B8, D1-D3 and E1-E2 hold. Let the marginal regressors  $X_s$  satisfy C1, and take  $\tau = \tau_n \sim (\log n)$ . Then the estimator  $\hat{m}_\tau(\cdot)$  with respect to sample observations  $\{Y_t, X_t\}_{t=1}^n$  satisfying A1' is uniformly consistent for  $m(x) = m(x_1, x_2, \dots)$  over  $S_\tau$ :*

$$\sup_{x \in S_\tau} |\hat{m}_\tau(x) - m_\tau(x)| = O_P \left( h^\beta + \sqrt{\frac{(\log n)^2 \exp(\kappa h^{-2/(2p-1)})}{nh^{\frac{1-p}{2p-1}}}} \right). \quad (45)$$

*If alternatively  $X_s$  is Gaussian and satisfies C2, then the same conclusion holds with respect to sample observations satisfying either A1' or A2'.*

**REMARK.** We may choose the optimal bandwidth as before; following the same arguments in the pointwise case, choosing  $h \sim (\log n)^a$  and solving for  $n$  gives

$$a_{opt} = \frac{\vartheta \cdot \mathcal{W} \left[ \frac{\chi}{\vartheta} c \exp(-\frac{\chi}{\vartheta} 2 \log \log n + \chi \log n) \right] + 2\chi \log \log n - \chi \log n}{\vartheta \chi \log \log n}. \quad (46)$$

And because the order of the leading terms is  $(\log n)^{-(2p-1)/2}$  as in the pointwise case, it is straightforward to see that the lower bound of the optimal bandwidth in Corollary 2 still continues to hold; that is,  $h_{opt} \succeq (\log n)^{-(2p-1)/2}$ . This is again invariant to the choice of distribution  $F$  of the squared regressor. It is important to note that as before, potential cross-sectional dependence between marginal regressors and also their distributional properties are represented via  $c$ , the collection of constants that appear inside the exponential terms either in (32) and (33).

In the following lemma we derive the lower rate of convergence.

**LEMMA 3.** *The lower bound for estimating the minimax quadratic risk of the regression function  $m(\cdot)$  at a fixed point  $x$  is  $(\log n)^{-(2p-1)}$ ; specifically,*

$$\inf_{\tilde{m}} \sup_{m \in \mathcal{M}_\beta} E|\tilde{m}(x_0) - m(x_0)|^2 \succeq h_{opt}^{2\beta} \left( \succeq (\log n)^{-(2p-1)} \right) \quad (47)$$

where  $\mathcal{M}_\beta$  is the class of regression functions  $m$  satisfying Assumption B7, and  $h_{opt}$  is the optimal bandwidth that balances the squared bias and the variance, provided that there exists some constant  $c$  such that for any real number  $\theta$  we have

$$\int \left\{ \sqrt{p_\varepsilon(t)} - \sqrt{p_\varepsilon(t + \theta)} \right\}^2 dt \leq c\theta^2, \quad (48)$$

where  $p_\varepsilon(\cdot)$  is the density of the error  $\varepsilon$  defined in (2).

The results altogether give the optimal rate of convergence of our estimator as follows. The same argument trivially applies to the pointwise case.

**COROLLARY 4.** *Suppose conditions assumed in Theorem 4 and (47) hold. Upon choosing  $h \sim (\log n)^{a_{opt}}$ , where  $a_{opt}$  is as defined in (46), we have*

$$\sup_{x \in S_\tau} \left| \hat{m}_\tau(x) - m_\tau(x) \right| = O_P \left( [\log n]^{\beta \cdot a_{opt}} \right), \quad (49)$$

and this rate of convergence is minimax optimal in view of (47).

## 5 Application to the Risk Return Relationship

The relation between the expected excess return on the aggregate stock market - the so called "equity risk premium" - and its conditional variance has long been the subject of both theoretical and empirical research in financial economics. The risk-return relation is an important ingredient in optimal portfolio choice, and is central to the development of theoretical asset-pricing models aimed at explaining a host of observed stock market patterns. Asset pricing models generally predict a positive relationship between the risk premium on the market portfolio and the variance of its return. In

an influential paper, Merton (1973) obtained very simple restrictions albeit under somewhat drastic assumptions; he showed in the context of a continuous time partial equilibrium model that

$$\mu_t = E[(r_{mt} - r_{ft})|\mathcal{F}_{t-1}] = \gamma \text{var}[(r_{mt} - r_{ft})|\mathcal{F}_{t-1}] = \gamma \sigma_t^2, \quad (50)$$

where  $r_{mt}$ ,  $r_{ft}$  are the returns on the market portfolio and risk-free asset respectively, while  $\mathcal{F}_{t-1}$  is the market wide information available at time  $t - 1$ . The positive constant  $\gamma$  is the Arrow–Pratt measure of relative risk aversion. The linear functional form actually only holds when  $\sigma_t^2$  is constant; otherwise  $\mu_t$  and  $\sigma_t^2$  can be nonlinearly related, Gennotte and Marsh (1993). Further examples with a positive risk return trade-off include the external habit model of Campbell and Cochrane (1999) and the Long Run Risks model of Bansal and Yaron (2004). However, a negative risk-return relation is not inconsistent with (a general enough) equilibrium, Backus and Gregory (1993). Unfortunately, the empirical evidence on the risk-return relation is mixed and inconclusive. Ghysels, Santa-Clara, and Valkanov (2005), Lundblad (2005), Pástor, Sinha, and Swaminathan (2008), and Ludvigson and Ng (2007) find a positive risk-return relation, while Campbell (1987), Glosten, Jagannathan, and Runkle (1993), Harvey (2001), and Lettau and Ludvigson (2003) find a negative relation. Still others find mixed and inconclusive evidence like French, Schwert, and Stambaugh (1987), Nelson (1991), Campbell and Hentschel (1992), Linton and Perron (2003), and Whitelaw (1994). Scruggs (1998) and Guo and Whitelaw (2006) document a positive trade-off within specifications that facilitate hedging demands. However, Scruggs and Glabadanidis (2003) find that this partial relationship is not robust across alternative volatility specifications. The main difficulty in estimating the risk-return relation is that neither the conditional expected return nor the conditional variance of the market is directly observable. The contradictory findings of the above studies are mostly the result of differences in the approaches to modeling the conditional mean and variance. Some studies have relied on parametric and semi-parametric ARCH or stochastic volatility models that impose a high degree of structure on the return generating process, about which there is little direct empirical evidence. Other studies have typically measured the conditional expectations underlying the conditional mean and conditional variance as projections onto predetermined conditioning variables. Practical constraints, such as choosing among a few conditioning variables, introduce an element of arbitrariness into the econometric modeling of expectations and can lead to omitted information estimation bias.

Pagan and Hong (1990) initiated the use of nonparametric methods in this setting. They argued that the risk premium  $\mu_t$  and the conditional variance  $\sigma_t^2$  are highly nonlinear functions of the past whose form is not captured by standard parametric GARCH–M models. They estimated  $E(r_{mt} - r_{ft}|r_{m,t-1}, \dots, r_{m,t-p})$  and  $\text{var}(r_{mt} - r_{ft}|r_{m,t-1}, \dots, r_{m,t-p})$  nonparametrically, where  $p \in \{1, 5\}$ , finding evidence of considerable nonlinearity. They then estimated  $\delta$  from the regression

$$r_{mt} - r_{ft} = \delta \sigma_t^2 + \eta_t, \quad (51)$$

by OLS and IV methods, finding a negative but insignificant  $\delta$ . There are a number of drawbacks with the Pagan and Hong (1990) approach. Firstly, the conditional moments are calculated using a

finite conditioning set. This greatly restricts the dynamics for the variance process. Secondly, they only test for linearity of the relationship between  $\mu_t$  and  $\sigma_t^2$ ; this seems to be somewhat restrictive in view of earlier findings. Linton and Perron (2003) considered the model where  $\sigma_t^2$  was a parametrically specified CH process (with dependence on the infinite past) but  $\mu_t = \varphi(\sigma_t^2)$  for some function  $\varphi$  of unknown functional form. They proposed an estimation algorithm but did not establish any statistical properties. They found some evidence of a nonlinear relationship. Conrad and Mammen (2008) develop the theory of estimation and inference for this model. Christensen, Dahl, and Iglesias (2012) developed the theoretical framework by considering volatility models that are driven by observable shocks so that a full theory can be given. Escanciano, Pardo-Fernández and Van Keilegom (2015) consider a more general class of semiparametric models.

Under the semi-strong form of the efficient markets hypothesis prices contain all relevant information and so the risk premium and risk themselves can be expressed in terms of only the past history of prices. We shall use this assumption to obviate the omitted variables/endogeneity issues that have limited previous applications in this area. Let  $\mu(x) = E(Y|X = x)$  and  $\sigma^2(x) = \text{var}(Y|X = x)$ , where  $Y$  is aggregate stock market returns in excess of the risk free rate and  $X$  is lagged values of returns. We suppose that both functions are unrestricted nonparametric functions of the entire information set and they are related in a quadratic way, that is,

$$\mu(x) = \alpha + \beta\sigma(x) + \gamma\sigma^2(x),$$

where  $\theta = (\alpha, \beta, \gamma)^\top$  with  $\alpha, \beta, \gamma$  being unknown constants. Let  $x_1, x_2, \dots, x_q \in \mathbb{R}^\infty$  be some given points such that  $\|D^{-1}(x_j - x_k)\| > 0$  for all  $j, k$ , and let  $\hat{\mu}(x)$  and  $\hat{\sigma}^2(x)$  be the estimated moments. Then we take

$$\hat{\theta} = (\hat{\alpha}, \hat{\beta}, \hat{\gamma})^\top = \hat{\Sigma}_q^{-1} \hat{U}_q$$

$$\hat{\Sigma}_q = \begin{pmatrix} 1 & \sum_{i=1}^q \hat{\sigma}(x_i) & \sum_{i=1}^q \hat{\sigma}^2(x_i) \\ \sum_{i=1}^q \hat{\sigma}(x_i) & \sum_{i=1}^q \hat{\sigma}^2(x_i) & \sum_{i=1}^q \hat{\sigma}^3(x_i) \\ \sum_{i=1}^q \hat{\sigma}^2(x_i) & \sum_{i=1}^q \hat{\sigma}^3(x_i) & \sum_{i=1}^q \hat{\sigma}^4(x_i) \end{pmatrix} \quad ; \quad \hat{U}_q = \begin{pmatrix} \sum_{i=1}^q \hat{\mu}(x_i) \\ \sum_{i=1}^q \hat{\sigma}(x_i) \hat{\mu}(x_i) \\ \sum_{i=1}^q \hat{\sigma}^2(x_i) \hat{\mu}(x_i) \end{pmatrix},$$

where  $q$  is finite.

We next derive the limiting distribution of the vector of estimated coefficients  $\hat{\theta} := (\hat{\alpha}, \hat{\beta}, \hat{\gamma})^\top$ , which can be used for conducting statistical inference. Define:

$$\Sigma_q = \begin{pmatrix} 1 & \sum_{i=1}^q \sigma(x_i) & \sum_{i=1}^q \sigma^2(x_i) \\ \sum_{i=1}^q \sigma(x_i) & \sum_{i=1}^q \sigma^2(x_i) & \sum_{i=1}^q \sigma^3(x_i) \\ \sum_{i=1}^q \sigma^2(x_i) & \sum_{i=1}^q \sigma^3(x_i) & \sum_{i=1}^q \sigma^4(x_i) \end{pmatrix}$$

$$\Omega(x_i) = \begin{pmatrix} \sigma^2(x_i) & \text{skew}(Y_t|X_t = x_i) \\ \text{skew}(Y_t|X_t = x_i) & \sigma^4(x_i) (\text{kurt}(Y_t|X_t = x_i) + 2) \end{pmatrix} =: \begin{pmatrix} \omega_{1,1}(x_i) & \omega_{1,2}(x_i) \\ \omega_{2,1}(x_i) & \omega_{2,2}(x_i) \end{pmatrix}, \quad (52)$$

$$V_q = \sum_{i=1}^q J(x_i) \Omega(x_i) J(x_i)^\top \quad ; \quad J(x_i) = \begin{pmatrix} 1 & 0 \\ \sigma(x_i) & \frac{\mu}{2\sigma}(x_i) \\ \sigma^2(x_i) & \mu(x_i) \end{pmatrix}.$$

Here, skew and kurt denote skewness and kurtosis of  $Y_t$  (conditional on  $X_t = x_i$ ). The result is a direct consequence of consistency of estimated moments and their asymptotic independence across  $i$ .

**THEOREM 5.** *Let Assumptions B2, B3, B5-B8, and D1-D4 hold, and suppose B4 is strengthened to require  $E(|Y_t|^{8+\delta}) \leq C < \infty$  for some  $C, \delta > 0$ . Suppose the operator  $g(\cdot) = E(Y^2|X = \cdot)$  satisfies Assumption B7. Suppose further that  $\omega_{a,b}(u)$  is continuous in some neighbourhood of  $x_i$  for all  $i$ . Then, given the sample observations  $\{Y_t, X_t\}_{t=1}^n$  specified in A2, we have the following limiting distribution:*

$$\sqrt{nh^{\frac{1-p}{2p-1}} \exp\left(-\kappa'_0 h^{-\frac{2}{2p-1}}\right)} \left(\hat{\theta} - \theta - B_\theta\right) \Rightarrow N\left(0, \kappa_2(K, p, a) \Sigma_q^{-1} V_q \Sigma_q^{-1}\right),$$

where  $B_\theta$  is a bias terms of order  $h^\beta$ ,  $\kappa_2$  is the constant in (33).

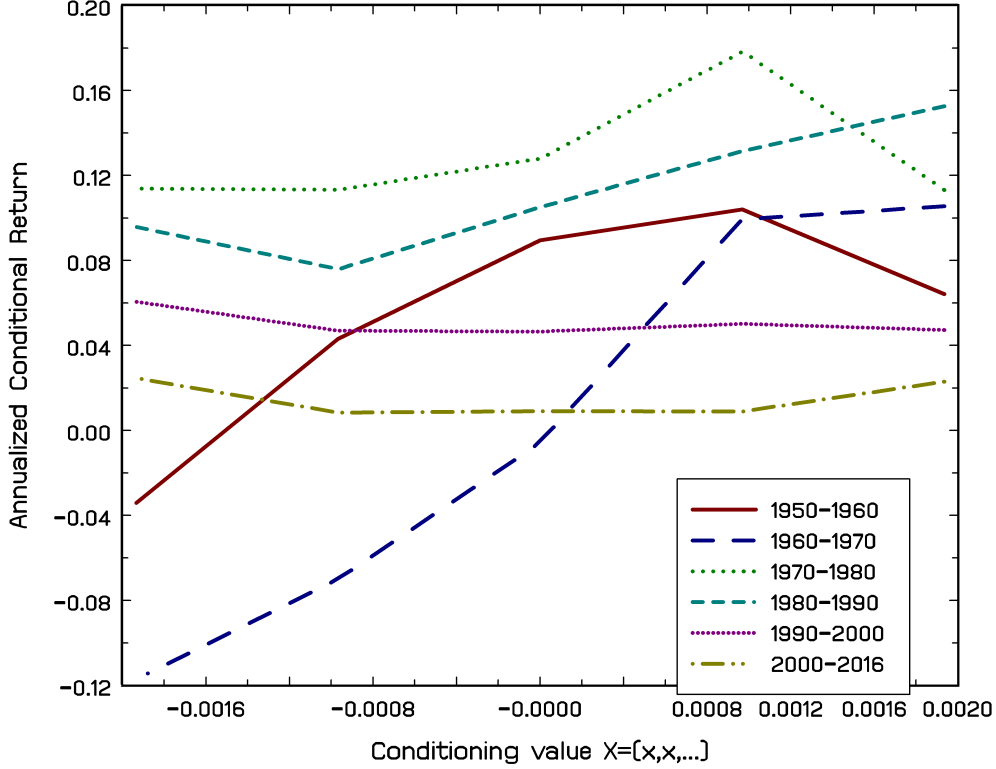
The parameters  $\theta$  are estimated at the same rate as the functions  $\mu(\cdot)$  and  $\sigma^2(\cdot)$ . It may be possible to achieve faster rates of convergence by allowing  $q \rightarrow \infty$ , as is commonly done in the semiparametric literature, but we have not yet been able to establish this rate improvement; see Chen and Christensen (2015).

We apply these methods to the daily risk premium on the value weighted S&P500 index — the total return on the index minus the returns on T-bills — over the period January 1950 to November 2016, a total of 16,820 observations. There is quite considerable variation in the ex-post risk and return by decades. We estimate  $\mu(x)$  and  $\sigma^2(x)$  using uniform kernel and the bandwidth sequence  $h = 0.5 \times s$ , where  $s$  was the full sample standard deviation, and  $p = 4$ . We estimated at the points  $X = (x, x, \dots)$  with  $x = 0, \pm 0.1s, \pm 0.2s$ .

**Table 1**

	Full and Sub Period						
	1950-2016	1950-1960	1960-1970	1970-1980	1980-1990	1990-2000	2000-2016
$255 \times \mu$	8.52	13.67	5.02	2.21	13.85	15.10	4.11
$\sqrt{255} \times \sigma^2$	15.45	11.53	10.25	13.69	17.23	14.18	19.94
$\alpha$	63.489	4.280	26.485	154.077	-284.338	29.075	1052.215
$\beta$	-3829.669	-275.651	-2399.650	-9733.012	14907.935	-1764.49	-51696.121
$\gamma$	57767.094	4275.594	54082.905	153819.71	-195301.47	26813.274	634971.52

We give the fitted surface by decade against the conditioning value. This suggests that the conditional risk premium was negative in the 1950s and 1960s when the conditioning value was strongly negative, whereas this phenomenon has disappeared in later decades with a flatter prediction profile.



## 6 Concluding remarks

Other quantities of interest in prediction such as the conditional median or mode can also be studied. This could be done via nonparametrically estimating the conditional distribution  $P(Y \leq y|X = \cdot) = E(1\{-\infty, y](Y)|X = \cdot)$ , but would necessarily require a slightly different set of assumptions. It is also quite easy to bring finite dimensional predictors into the theory separately. For example, one may want to allow for slow time variation whereby  $t/T$  becomes an additional covariate and the regression function is  $m(x, u)$  with  $u \in [0, 1]$  and  $x \in \mathbb{R}^\infty$ . In this case we modify the estimator of (17) by introducing a multiplicative kernel of the form  $k_b(u - t/T)$ , where  $b$  is a bandwidth and  $k$  is a symmetric probability density function.

## 7 Appendix: Proofs of the main results

### 7.1 Proof of Theorem 1

From the decomposition (20):

$$\hat{m}(x) - m(x) = \frac{E\hat{m}_2(x) - m(x)}{\hat{m}_1(x)} + \frac{\hat{m}_2(x) - E\hat{m}_2(x)}{\hat{m}_1(x)} - \frac{m(x)[\hat{m}_1(x) - 1]}{\hat{m}_1(x)},$$



we see that it suffices to show  $E\hat{m}_2(x) - m(x) \rightarrow 0$  and  $\hat{m}_2(x) - E\hat{m}_2(x) \xrightarrow{P} 0$ , since  $\hat{m}_1(x) \xrightarrow{P} 1$  would then follow from the latter and complete the proof.

As for the former ‘bias component’, denoting by  $\mathcal{E}(x, \lambda \underline{h})$  the infinite dimensional hyperellipsoid centred at  $x = (x_j)_j \in \mathbb{R}^\infty$  with semi-axes  $h_j$  in each direction we have

$$\begin{aligned} E\hat{m}_2(x) - m(x) &= E\left(\frac{1}{nEK_1} \sum_{t=1}^n K_t Y_t - m(x)\right) \\ &= \frac{1}{EK_1} EK_1 Y_1 - \frac{EK_1}{EK_1} m(x) = \frac{1}{EK_1} E\left[E\left[(Y_1 - m(x))K_1 \middle| X\right]\right] \\ &= \frac{1}{EK_1} E\left[\left[m(X) - m(x)\right]K_1\right] \leq \sup_{u \in \mathcal{E}(x, \lambda \underline{h})} |m(u) - m(x)| \longrightarrow 0 \end{aligned} \quad (53)$$

as  $n \rightarrow \infty$ , where  $K_t$  is the shorthand notation for  $K(\|H^{-1}(x - X_t)\|)$  as introduced in the main text before. The second equality is justified by stationarity that is preserved under measurable transformation, and the last inequality is due to compact support of the kernel and continuity of the regression operator at  $x$  (Assumption B1).

The next step concerns with the latter ‘variance component’  $\hat{m}_2 - E\hat{m}_2$ ; its mean-squared convergence to zero will be shown. Writing

$$\hat{m}_2 - E\hat{m}_2 = \frac{1}{n} \sum_{t=1}^n \frac{1}{EK_1} \left\{ K_t Y_t - E(K_t Y_t) \right\} =: \frac{1}{n} \sum_{t=1}^n Q_{nt}, \quad (54)$$

we remark that the arguments to follow depend upon the temporal dependence structure of  $Q_{nt}$ . In the static regression case,  $Q_{nt}$  is a measurable function of  $Y_t, X_{1t}, X_{2t}, \dots$ , and hence inherits their joint dependence structure. That is,  $Q_{nt}$  is arithmetically  $\alpha$ -mixing with the rate specified in A1. In the dynamic regressions case (which covers the autoregression framework), the dependence of  $Q_{nt}$  is defined via  $K_t$  which is near epoch dependent on  $(Y_t, V_t)$  as specified in Assumption A2; this bypasses the issue of  $Q_{nt}$  being dependent upon infinite past of  $Y_t$  and/or  $V_t$ . We proceed with these two cases separately.

**CASE 1: STATIC REGRESSION.** Clearly, it is sufficient to prove  $\text{var}(\hat{m}_2 - E\hat{m}_2) \rightarrow 0$  for mean squared convergence. Since  $Q_{nt}$  is stationary over time we have

$$\text{var}(\hat{m}_2 - E\hat{m}_2) = \frac{1}{n^2} \sum_{t=1}^n \text{var}(Q_{nt}) + \frac{2}{n^2} \sum_{1 \leq i < j \leq n} \text{cov}(Q_{ni}, Q_{nj}) \quad (55)$$

$$\begin{aligned} &= \frac{1}{n} \text{var}(Q_{n1}) + \frac{2}{n^2} \sum_{1 \leq j-i < n} \text{cov}(Q_{ni}, Q_{nj}) \\ &= \frac{1}{n} \text{var}(Q_{n1}) + \frac{2}{n^2} \sum_{s=1}^{n-1} (n-s) \cdot \text{cov}(Q_{n1}, Q_{n,s+1}) =: A_1 + A_2. \end{aligned} \quad (56)$$

Now, by (8), (10) and Assumption A it follows that

$$\begin{aligned} A_1 &= \frac{1}{nE^2K_1} \text{var} \left( K_1 Y_1 - EY_1 K_1 \right) = \frac{\text{var}(K_1 Y_1)}{nE^2K_1} \\ &\leq \frac{EK_1^2 Y_1^2}{nE^2K_1} = \frac{E(E(Y_1^2|X_1)K_1^2)}{nE^2K_1} \leq \frac{C}{n\varphi_x(\lambda \underline{h})} \rightarrow 0 \end{aligned} \quad (57)$$

as  $n \rightarrow \infty$ .

We now move on to the second term  $A_2$  and investigate the covariance term. Since measurable transformations of mixing variables preserve the mixing property, using Davydov's inequality, see Davydov (1968, Lemma 2.1) or Bosq (1996, Corollary 1.1) and stationarity we have

$$|\text{cov}(Q_{n1}, Q_{n,s+1})| = \left| \text{cov} \left( Y_1 \frac{K_1}{EK_1}, Y_{s+1} \frac{K_{s+1}}{EK_1} \right) \right| \leq \frac{C \{E|Y_1 K_1|^{2+\delta}\}^{\frac{2}{2+\delta}}}{\varphi_x(\underline{h}\lambda)^2 \cdot s^{k\delta/(2+\delta)}}. \quad (58)$$

In the meantime,

$$\begin{aligned} |\text{cov}(Q_{n1}, Q_{n,s+1})| &= \left| \text{cov} \left( Y_1 \frac{K_1}{EK_1}, Y_{s+1} \frac{K_{s+1}}{EK_1} \right) \right| \\ &\leq \left| E \left( Y_1 \frac{K_1}{EK_1} Y_{s+1} \frac{K_{s+1}}{EK_1} \right) \right| + \left| E \left( Y_1 \frac{K_1}{EK_1} \right) E \left( Y_{s+1} \frac{K_{s+1}}{EK_1} \right) \right| \\ &\leq \frac{C}{\varphi_x(\underline{h}\lambda)^2} |E(K_1 K_{s+1})| + \frac{C'}{E^2 K_1} |E(K_1) E(K_{s+1})| \\ &\leq \frac{C}{\varphi_x(\underline{h}\lambda)^2} \cdot \psi_x(\lambda \underline{h}; 1, s+1) + C' \leq C'' \end{aligned} \quad (59)$$

by stationarity, law of iterated expectation, boundedness of regression function, and Assumption B6, B5 (along with the upper bound  $\psi(\lambda \underline{h}; 1, s+1)$  of  $EK_1 K_{s+1}$  obtained as a direct consequence of B5 following similar arguments used for Lemma 1).

With reference to (58) and (59), we take some increasing sequence  $u_n \rightarrow \infty$  such that  $u_n = o(n)$ , and write

$$\begin{aligned} \sum_{s=1}^{n-1} |\text{cov}(Q_{n1}, Q_{n,s+1})| &= \sum_{s=1}^{u_n-1} |\text{cov}(Q_{n1}, Q_{n,s+1})| + \sum_{s=u_n}^{n-1} |\text{cov}(Q_{n1}, Q_{n,s+1})| \\ &\leq C''(u_n - 1) + \sum_{s=u_n}^{n-1} \frac{C s^{-k\delta/(2+\delta)}}{\varphi_x(\underline{h}\lambda)^2} = O \left( u_n + \frac{u_n^{-k\delta/(2+\delta)+1}}{\varphi_x(\underline{h}\lambda)^2} \right), \end{aligned} \quad (60)$$

which is  $O(\varphi_x(\underline{h}\lambda)^{-2(2+\delta)/(k\delta)})$  upon choosing  $u_n \sim \varphi_x(\underline{h}\lambda)^{-2(2+\delta)/(k\delta)}$ .

Consequently, since  $k \geq 2(2+\delta)/\delta$  it follows that

$$\begin{aligned} A_2 &:= \frac{2}{n^2} \sum_{s=1}^{n-1} (n-s) \cdot \text{cov}(Q_{n1}, Q_{n,s+1}) = \frac{2}{n} \sum_{s=1}^{n-1} \left( 1 - \frac{s}{n} \right) \cdot \text{cov}(Q_{n1}, Q_{n,s+1}) \\ &= O(n^{-1} [\varphi_x(\underline{h}\lambda)]^{-2(2+\delta)/(k\delta)} + n^{-2} [\varphi_x(\underline{h}\lambda)]^{-2(2+\delta)/(k\delta)}) \\ &= O(n^{-1} [\varphi_x(\underline{h}\lambda)]^{-2(2+\delta)/(k\delta)}) = o(1) \end{aligned} \quad (61)$$

by Assumption B2, and the desired result is obtained.

CASE 2: DYNAMIC REGRESSION.<sup>5</sup> We return back to (54):

$$\widehat{m}_2 - E\widehat{m}_2 = \frac{1}{n} \sum_{t=1}^n \frac{1}{EK_1} \left\{ K_t Y_t - E(K_t Y_t) \right\} =: \frac{1}{n} \sum_{t=1}^n Q_{nt}. \quad (62)$$

In this framework  $K_t = K(\|H^{-1}(x - X_t)\|)$  is a (measurable) function of  $(Y_{t-1}, Y_{t-2}, \dots)$ . Despite loosing the mixing property,  $K_t$  inherits stationarity of the mixing process  $\{Y_t\}$ . We write  $K_{t,(r)} = \Psi(Y_t, Y_{t-1}, Y_{t-2}, \dots, Y_{t-r+1}) = E(K_t | Y_t, \dots, Y_{t-r+1})$ , where  $\Psi$  denotes a measurable map and  $r$  is as in Assumption A2. Clearly,  $K_{t,(r)}$  preserves the mixing dependence structure of  $Y_t$  with mixing coefficient  $\alpha(\ell - (r - 1))$  since  $\sigma(K_{s,(r)}; s \geq t + \ell) \subset \sigma((Y_s, \dots, Y_{s-r+1}); s \geq t + \ell) = \sigma(Y_s; s \geq t + \ell - (r - 1))$ .

Now write

$$\begin{aligned} \widehat{m}_2 - E\widehat{m}_2 &= \frac{1}{n} \sum_{t=1}^n \frac{1}{EK_1} \left[ K_{t,(r)} Y_t - E(K_{t,(r)} Y_t) \right] + \frac{1}{n} \sum_{t=1}^n \frac{1}{EK_1} \left[ K_t Y_t - K_{t,(r)} Y_t \right] \\ &\quad + \frac{1}{n} \sum_{t=1}^n \frac{1}{EK_1} \left[ E(K_{t,(r)} Y_t) - E(K_t Y_t) \right] = R_1 + R_2 + R_3, \end{aligned} \quad (63)$$

and first consider the last term  $R_3$ .

Fix some increasing sequence  $q = q_n \rightarrow \infty$ , and write  $Y_{t,L} := Y_t 1\{|Y_t| \leq q\}$  and  $Y_{t,U} = Y_t 1\{|Y_t| > q\}$ . Then

$$\begin{aligned} EY_t K_{t,(r)} &= EY_t K(\|H^{-1}(x - X_t)\|) - EY_{t,U} K(\|H^{-1}(x - X_t)\|) \\ &\quad + EY_{t,L} K_{t,(r)} - EY_{t,L} K(\|H^{-1}(x - X_t)\|) \\ &\quad + EY_{t,U} K_{t,(r)} = D_1 + D_2 + D_3. \end{aligned} \quad (64)$$

The second part of  $D_1$  is given by

$$\begin{aligned} EY_{t,U} K(\|H^{-1}(x - X_t)\|) &\leq E|Y_t| 1_{\{|Y_t| > q\}} K(\|H^{-1}(x - X_t)\|) \\ &\leq q^{-(\delta+1)} E|Y_t|^{2+\delta} 1_{\{|Y_t| > q\}} K_t \leq Cq^{-(\delta+1)} E|Y_t|^{2+\delta} 1_{\{|Y_t| > q\}} = o(q^{-(\delta+1)}) \end{aligned} \quad (65)$$

because  $1_{\{|Y_t| > q\}} = o(1)$  as  $n \rightarrow \infty$ . Following similar arguments on  $D_3$  we have  $D_1 + D_3 = EY_t K_t + o(q^{-(\delta+1)})$ . So we are now left with the middle term  $D_2$ :

$$D_2 \leq E|Y_{t,L}| |K_t - K_{t,(r)}| = O\left(q\sqrt{v_2(r_n)}\right) \quad (66)$$

by Hölder's inequality. Therefore, from (64), (65) and (66) we see that

$$R_3 = \frac{1}{nEK_1} \sum_{t=1}^n \left[ EK_{t,(r)} Y_t - E(K_t Y_t) \right] = o\left(\frac{q^{-(\delta+1)}}{\varphi_x(\lambda \underline{h})}\right) + O\left(\frac{q\sqrt{v_2(r_n)}}{\varphi_x(\lambda \underline{h})}\right), \quad (67)$$

---

<sup>5</sup>For the sake of notational simplicity, we will write the proofs for the dynamic regression framework in terms of its autoregressive special case throughout the appendix. That is, some lags of the response variable  $Y_t$  here possibly represent the lags of the covariate  $V_t$ .

and upon choosing  $q = (\varphi_x(\underline{h}\lambda)/n)^{-1/(2(\delta+1))}$  we have  $o(\varphi^{-1}q^{-(\delta+1)}) = o(\varphi^{-1}(\varphi/n)^{1/2}) = o(n^{-1/2}\varphi^{-1/2}) = o(1)$ . Furthermore,

$$\begin{aligned} O\left(\frac{1}{\varphi_x(\underline{h}\lambda)}q\sqrt{v_2(r_n)}\right) &= O\left(\frac{1}{\varphi_x(\underline{h}\lambda)} \cdot \left(\frac{\varphi_x(\underline{h}\lambda)}{n}\right)^{-1/(2(\delta+1))} \sqrt{v_2(r_n)}\right) \\ &= O\left(\frac{\sqrt{v_2(r_n)}}{[\varphi_x(\underline{h}\lambda)]^{(2\delta+3)/(2\delta+2)}n^{-1/(2(\delta+1))}}\right) = o(1) \end{aligned} \quad (68)$$

by Assumption A2, yielding  $R_3 = o(1)$ , and hence  $R_2 = o_p(1)$ .

As for the first term that remains,

$$\begin{aligned} R_1 &= \frac{1}{n} \sum_{t=1}^n \left[ \frac{K_{t,(r)}Y_t - E(K_tY_t)}{EK_1} \right] + \frac{1}{n} \sum_{t=1}^n \left[ \frac{E(K_tY_t) - E(K_{t,(r)}Y_t)}{EK_1} \right] \\ &= \frac{1}{n} \sum_{t=1}^n E(Q_{nt}|Y_t, Y_{t-1}, \dots, Y_{t-r+1}) - R_3 \\ &= \frac{1}{n} \sum_{t=1}^n Q_{nt,(r)} + o\left(\frac{q^{-(\delta+1)}}{\varphi_x(\underline{h}\lambda)}\right) + O\left(\frac{\sqrt{v_2(r_n)}}{[\varphi_x(\underline{h}\lambda)]^{(2\delta+3)/(2\delta+2)}n^{-1/(2(\delta+1))}}\right). \end{aligned} \quad (69)$$

Since  $Q_{nt,(r)}$  is  $\alpha$ -mixing, we can work with the first term by following similar arguments in the regression case. Specifically, due to boundedness of the kernel and the mixing properties, the bound in (58) can be constructed. As for the constant bound constructed in (59), we rewrite

$$\begin{aligned} \frac{\text{cov}(Y_1K_{1,(r)}, Y_{s+1}K_{s+1,(r)})}{\varphi_x(\lambda\underline{h})^2} &= \frac{\text{cov}(Y_1[K_{1,(r)} - K_1], Y_{s+1}[K_{s+1,(r)} - K_{s+1}])}{\varphi_x(\lambda\underline{h})^2} \\ &\quad + \frac{\text{cov}(Y_1[K_{1,(r)} - K_1], Y_{s+1}K_{s+1,(r)})}{\varphi_x(\lambda\underline{h})^2} \\ &\quad + \frac{\text{cov}(Y_1, Y_{s+1}[K_{s+1,(r)} - K_{s+1}])}{\varphi_x(\lambda\underline{h})^2} + \frac{\text{cov}(Y_1K_1, Y_{s+1}K_{s+1})}{\varphi_x(\lambda\underline{h})^2} \\ &= \mathcal{G}_1 + \mathcal{G}_2 + \mathcal{G}_3 + \mathcal{G}_4. \end{aligned}$$

By (58),  $\mathcal{G}_4 \leq C$  by (59). Further,

$$\begin{aligned} \mathcal{G}_1 &\leq \left| \frac{E(Y_1Y_{s+1}[K_{1,(r)} - K_1][K_{s+1,(r)} - K_{s+1}])}{\varphi_x(\lambda\underline{h})^2} \right| \\ &\quad + \left| \frac{E(Y_1[K_{1,(r)} - K_1]) \cdot E(Y_{s+1}[K_{s+1,(r)} - K_{s+1}])}{\varphi_x(\lambda\underline{h})^2} \right| \leq C' \frac{v_2(r)}{\varphi_x(\lambda\underline{h})^2} \rightarrow 0 \end{aligned}$$

by Assumption B6 and by the fact that

$$\left( \frac{\sqrt{v_2(r_n)}}{\varphi_x(\underline{h}\lambda)} \right) \leq \left( \frac{\sqrt{v_2(r_n)}}{\varphi_x(\underline{h}\lambda)} \right) \cdot (n/\varphi)^{1/(2\delta+2)} \rightarrow 0$$

by (18) in Assumption A2. Similarly,  $\mathcal{G}_2$  and  $\mathcal{G}_3$  can be easily shown to converge to zero in large sample.

Now choosing an increasing sequence  $u_n \sim [\varphi_x(\underline{h}\lambda)^{-2(2+\delta)/(k\delta)} + r_n] \rightarrow \infty$  such that  $r_n/u_n = o(1)$ , we see that (ignoring the array notation in  $Q_{nt,(r)}$  for simplicity)

$$\begin{aligned} \sum_{s=1}^{n-1} |\text{cov}(Q_{1,(r)}, Q_{s+1,(r)})| &= \sum_{s=1}^{u_n-1} |\text{cov}(Q_{1,(r)}, Q_{s+1,(r)})| + \sum_{s=u_n}^{n-1} |\text{cov}(Q_{1,(r)}, Q_{s+1,(r)})| \\ &\leq C(\varphi_x(\underline{h}\lambda)^{-\frac{2(2+\delta)}{(k\delta)}} + r_n) + \sum_{s=u_n}^{n-1} \frac{C(s - r_n + 1)^{-k\delta/(2+\delta)}}{\varphi_x(\underline{h}\lambda)^2} = O\left(\varphi_x(\underline{h}\lambda)^{-\frac{2(2+\delta)}{(k\delta)}}\right), \end{aligned}$$

since the mixing coefficient for  $Q_{nt,(r)}$  denoted  $\alpha'(n)$  is given by  $\alpha(n - (r - 1))$  for  $n \geq r$ . It now follows by the same arguments in (61) that the first term in (69) converges to zero, and hence  $R_1 = o_p(1)$ , which is the result we desired.  $\blacksquare$

## 7.2 Proof of Theorem 2 and 3

We start by recalling the bias component discussed in (53). Additional assumptions B7, B8 and D3 allow us to proceed further as follows:

$$\begin{aligned} \mathcal{B}_n(x) &= E\hat{m}_2(x) - m(x) = E\left(\frac{1}{nEK_1} \sum_{t=1}^n K_t Y_t - m(x)\right) \\ &= \frac{1}{EK_1} EK_1 Y_1 - \frac{EK_1}{EK_1} m(x) = \frac{1}{EK_1} E\left[E\left[(Y_1 - m(x))K_1 \middle| X\right]\right] \\ &= \frac{1}{EK_1} E\left[\left[m(X) - m(x)\right]K_1\right] \leq \sup_{u \in \mathcal{E}(x, \lambda \underline{h})} |m(u) - m(x)| \\ &\leq \sup_{u \in \mathcal{E}(x, \lambda \underline{h})} \sum_{j=1}^{\infty} c_j |u_j - x_j|^\beta = \sum_{j=1}^{\infty} c_j (\lambda h \phi_j)^\beta = h^\beta \left(\lambda^\beta \sum_{j=1}^{\infty} c_j j^{p\beta}\right) < \infty. \end{aligned} \quad (70)$$

Now rewriting the decomposition (20) as

$$\begin{aligned} \hat{m}(x) - m(x) - \mathcal{B}_n(x) &= \frac{\mathcal{B}_n(x) \cdot [1 - \hat{m}_1(x)]}{\hat{m}_1(x)} + \frac{\hat{m}_2(x) - E\hat{m}_2(x) - m(x)[\hat{m}_1(x) - 1]}{\hat{m}_1(x)}, \end{aligned}$$

and noting that  $\hat{m}_1(x) \rightarrow^p 1$  (an immediate consequence of Theorem 1), we see that it suffices to derive the limiting distribution of

$$\begin{aligned} \hat{m}_2(x) - E\hat{m}_2(x) - m(x)[\hat{m}_1(x) - 1] &= \frac{1}{n} \sum_{t=1}^n \frac{1}{EK_1} \left[ K_t Y_t - m(x) K_t - E(K_t Y_t) + m(x) EK_t \right] =: \frac{1}{n} \sum_{t=1}^n R_{nt}. \end{aligned} \quad (71)$$

By Assumption B6, D3, D4, and the law of iterated expectations, the asymptotic variance of the triangular array  $R_{nt}$  is given by

$$\begin{aligned}
\text{var}(R_{nt}) &= \frac{\text{var}[K_t(Y_t - m(x))]}{E^2 K_1} \\
&= \frac{1}{E^2 K_1} \left\{ E \left[ K_t(Y_t - m(x)) \right]^2 - E^2 \left[ K_t(Y_t - m(x)) \right] \right\} \\
&\simeq \frac{1}{E^2 K_1} \left\{ E \left[ \sigma^2(X) K_1^2 \right] + E \left( \left[ m(X) - m(x) \right]^2 K_1^2 \right) \right\} \\
&= \frac{1}{E^2 K_1} \left\{ \sigma^2(x) E K_1^2 + E \left( \left[ \sigma^2(X) - \sigma^2(x) \right] K_1^2 \right) + o(1) E K_1^2 \right\} \\
&= \frac{E K_1^2}{E^2 K_1} (\sigma^2(x) + o(1)) \simeq \frac{\sigma^2(x) \xi_2}{\varphi_x(h\lambda) \xi_1^2}. \tag{72}
\end{aligned}$$

Following similar arguments and using the latter assumption of D4, it can be readily shown that the covariance term is of smaller order than (72), which together shows (30). Under Assumption D1 the small ball probability can be written in terms of the centered small deviation and  $p^*(\cdot)$ , the Radon-Nikodym derivative of the induced probability measure  $P_{z-Z}$  with respect to  $P_Z$ :

$$\begin{aligned}
\varphi_x(\lambda h) &= P(X \in \mathcal{E}(x, \lambda h)) \\
&= P \left( \sum_{j=1}^{\infty} j^{-2p} (x_j - X_j)^2 \leq h^2 \lambda^2 \right) = P \left( \|z - Z\| \leq h\lambda \right) \\
&= \int_{B(0, h\lambda)} dP_{z-Z}(u) = \int_{B(0, h\lambda)} p^*(u) dP_Z(u) \\
&\simeq p^*(0) \cdot P(\|Z\| \leq h\lambda) = p^*(0) \times P \left( \sum_{j=1}^n j^{-2p} X_j^2 \leq h^2 \lambda^2 \right), \tag{73}
\end{aligned}$$

where the latter probability can be explicitly specified by substituting  $r = h^2 \lambda^2$ ,  $A = 2p$ , and  $a = 2p/(2p-1)$  in Proposition 4.1 of Dunker et al. (1998) for the i.i.d. case. As for the case where the regressors are dependent, i.e. when the  $X_j$ 's satisfy Assumption C2, the small ball probability can be specified in view of Theorem 1.1 of Hong, Lifshits and Nazarov (2016). Finally we have,

$$\frac{\sigma^2(x) \xi_2}{\varphi_x(h\lambda) \xi_1^2} = \frac{1}{\phi(h)} \cdot \frac{\sigma^2(x) \xi_2}{p^*(0) \xi_1^2} \cdot \frac{C^* C_\ell}{\lambda^{\frac{1+2\rho p}{2p-1}}},$$

where  $\phi(h) = h^{(1+2\rho p)/(2p-1)} \exp\{-C^{**}(\lambda h)^{-2/(2p-1)}\}$ , and (as defined before)

$$C_\ell = \lim_{h \rightarrow 0} \left[ \ell^{-1/2} \left( h^{-\frac{4p}{2p-1}} \right) \right] \quad C^* = \frac{(2\pi)^{(1+2\rho p)} (2p-1)}{\Gamma^{-1}(1-\rho) \cdot (2p)^{\frac{2p(\rho+2)-1}{2p-1}}} \cdot \zeta^{\frac{2p(1+\rho)}{2p-1}}$$

and  $\Gamma(\cdot)$  is the Gamma function,  $\xi_1$  and  $\xi_2$  are the constants specified in (11), and  $\lambda$  is the upper bound of the support of the kernel.

In constructing the central limit theorem we consider the normalized statistic  $R_{nt}^* := \sqrt{\phi(h)} \cdot R_{nt}$  and derive the self normalized limiting distribution of  $(1/\sqrt{n}) \cdot R_{nt}^*$ . We shall only prove the autoregression case, where an additional step of mixing approximation is added to the standard regression case; the asymptotic normality for the regression case in a functional context was established in Masry (2005). In many places of the remainder of this proof we shall closely follow their proof of Theorem 4.

We make use of Bernstein's blocking method and partition  $\{1, \dots, n\}$  by  $2k (= 2k_n \rightarrow \infty)$  number of blocks of two different sizes that alternate (hereafter referred to as the "big" and "small" blocks) and lastly a single block (the "last block") that covers the remainder. The size of the alternating blocks is given by  $a_n$  and  $b_n$  respectively, where the one for the "big-blocks"  $a_n$  is set to dominate that for the "small-blocks"  $b_n$  in large sample, i.e.  $b_n = o(a_n)$ . More specifically, we take

$$k_n = \lfloor n/(a_n + b_n) \rfloor \quad \text{and} \quad a_n = \lfloor \sqrt{n\phi(h)}/q_n \rfloor$$

where  $q_n \rightarrow \infty$  is a sequence of integer; it then clearly follows that  $a_n/n \rightarrow 0$  and  $a_n/\sqrt{n\phi(h)} \rightarrow 0$ . We also assume  $(n/a_n) \cdot \alpha^*(b_n) = (n/a_n) \cdot \alpha(b_n - r + 1) \rightarrow 0$ , where  $\alpha^*$  is the mixing coefficient of  $R_{nt,(r)}^* = E(R_{nt}^* | \mathcal{F}_{t-r+1}^{t-1})$ .

By construction above we can write  $\sqrt{n}^{-1} \sum_{t=1}^n R_{nt}^*$  as the sum of the groups of big-blocks  $\mathcal{B}$ , small-blocks  $\mathcal{S}$  and the remainder block  $\mathcal{R}$  defined as

$$\begin{aligned} \mathcal{B} &:= \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \Xi_{1,j} = \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \left( \sum_{t=j(a+b)+1}^{j(a+b)+a} R_{nt}^* \right) \\ \mathcal{S} &:= \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \Xi_{2,j} = \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \left( \sum_{t=j(a+b)+a+1}^{(j+1)(a+b)} R_{nt}^* \right) \\ \mathcal{R} &:= \frac{1}{\sqrt{n}} \Xi_{3,j} = \frac{1}{\sqrt{n}} \left( \sum_{t=k(a+b)+1}^n R_{nt}^* \right). \end{aligned}$$

The aim is to show that the contributions from the small and the last remaining block are negligible, and that the big-blocks are asymptotically independent.

We first consider the big blocks  $\mathcal{B}$ . Given  $r$  as in Assumption 2, and  $R_{nt,(r)}^* = E(R_{nt}^* | Y_t, \dots, Y_{t-r+1})$  we have

$$\mathcal{B} = \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \left( \sum_{t=j(a+b)+1}^{j(a+b)+a} R_{nt,(r)}^* \right) + \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \left( \sum_{t=j(a+b)+1}^{j(a+b)+a} [R_{nt,(r)}^* - R_{nt}^*] \right) = \mathcal{Q}_1 + \mathcal{Q}_2.$$

As for the second term, consider

$$\begin{aligned}
\frac{1}{\sqrt{n}} E \mathcal{Q}_2 &\leq \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \sum_{t=j(a+b)+1}^{j(a+b)+a} E |R_{nt,(r)}^* - R_{nt}^*| \\
&= \frac{1}{EK_1} \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \sum_{t=j(a+b)+1}^{j(a+b)+a} E |K_t Y_t - Y_t E(K_t | Y_t, Y_{t-1}, \dots, Y_{t-r+1})| \\
&\leq \frac{1}{\sqrt{n}} \frac{1}{\varphi_x(\underline{h}\lambda)} \sum_{j=0}^{k-1} \sum_{t=j(a+b)+1}^{j(a+b)+a} E |Y_t| |K_t - K_{t,(r)}| \\
&\leq \frac{1}{\sqrt{n}} \frac{1}{\varphi_x(\underline{h}\lambda)} \sum_{j=0}^{k-1} \sum_{t=j(a+b)+1}^{j(a+b)+a} \left( E |Y_t|^2 \right)^{1/2} \left( E |K_t - K_{t,(r)}|^2 \right)^{1/2} \\
&\leq C \cdot \frac{1}{\sqrt{n}} k_n a_n \frac{\sqrt{v_2(r_n)}}{\varphi_x(\lambda \underline{h})} = O \left( \frac{\sqrt{n \cdot v_2(r_n)}}{\varphi_x(\lambda \underline{h})} \right) = o(1),
\end{aligned}$$

which implies that  $\sqrt{n}^{-1} \mathcal{Q}_2 = o_p(1)$ .

We now show asymptotic independence of terms in  $\mathcal{Q}_1$ , on noting that  $\Xi'_{1,j}$ s are independent if for all real  $t_j$

$$\left| E \left[ \sum_{j=0}^{k-1} \exp(it_j \Xi_{1,j}) \right] - \prod_{j=0}^{k-1} E[\exp(it_j \Xi_{1,j})] \right| \quad (74)$$

is zero, see for instance Applebaum (2009, page 18). Applying the Volkonskii-Rozanov inequality (see Fan and Yao (2003, page 72)), it can be shown that (74) is bounded above by  $C(n/a_n) \cdot \alpha(b_n - r + 1) \rightarrow 0$ , implying asymptotic independence.

Moving on to the small blocks, due to stationarity we have

$$\begin{aligned}
\text{var}(\mathcal{S}) &= \frac{1}{n} \text{var} \left( \sum_{j=0}^{k-1} \sum_{t=j(a+b)+a+1}^{(j+1)(a+b)} R_{nt}^* \right) \\
&= \frac{1}{n} \sum_{j=0}^{k-1} \text{var} \left( \sum_{t=j(a+b)+a+1}^{(j+1)(a+b)} R_{nt}^* \right) + \frac{1}{n} \sum_{j \neq l}^{k-1} \text{cov} \left( \sum_{t=j(a+b)+a+1}^{(j+1)(a+b)} R_{nt}^*, \sum_{s=l(a+b)+a+1}^{(l+1)(a+b)} R_{ns}^* \right) \\
&= \frac{1}{n} \sum_{j=0}^{k-1} \left( b_n \text{var}(R_{nt}^*) + \sum_{t \neq l}^{b_n} \text{cov}(R_{nt}^*, R_{nl}^*) \right) + \frac{1}{n} \sum_{j \neq l}^{k-1} \sum_{i,j=1}^{b_n} \text{cov}(R_{n,i+w_j}^*, R_{n,r+w_l}^*) \\
&= Q_1 + Q_2 + Q_3.
\end{aligned}$$

where  $w_j = j(a+b) + a$ .

Regarding the first term, similar arguments used in deriving (72) yield

$$Q_1 = \frac{1}{n} k_n b_n \frac{[\varphi_x(\underline{h}\lambda)^{1/2}]^2 \sigma^2(x) \xi_2}{\varphi_x(\underline{h}\lambda) \xi_1^2} = \frac{k_n b_n \sigma^2(x) \xi_2}{n \xi_1^2} \rightarrow 0 \quad (75)$$



because  $k_n b_n/n \sim b_n/(a_n + b_n) \rightarrow 0$ . Now moving on to  $Q_2$  and  $Q_3$ , the sum of covariances can be dealt with in the same manner as we did for the variance using (72), so  $Q_2 \rightarrow 0$ . Similarly for  $Q_3$ , implying  $\text{var}(\mathcal{S}) \rightarrow 0$  as desired. Convergence result for the remainder  $\mathcal{R}$  can be established similarly, and is bounded by  $C(a_n + b_n)/n \rightarrow 0$ .

The results above suggest that

$$\frac{1}{\sqrt{n}} \sum_{t=1}^n R_{nt}^* = \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \left( \sum_{t=j(a+b)+1}^{j(a+b)+a} R_{nt}^* \right) + o_p(1) = \frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \eta_j + o_p(1), \quad (76)$$

and the desired result holds in view of (62) and the CLT for triangular array upon checking the Lindeberg condition (which is omitted here due to its similarity with Masry (2005, page 174-175)). Corollary 2 now follows because

$$\begin{aligned} \sqrt{n\phi(h)} \left( \frac{\hat{m} - m - \mathcal{B}_n}{\sqrt{n\phi(h)\Delta_n}} \right) &= \frac{\sqrt{n} \frac{1}{n} \sum_{t=1}^n R_{nt}^*}{\sqrt{\frac{1}{n} \sum_t \hat{R}_{nt}^{*,2}}} = \frac{\frac{1}{\sqrt{n}} \sum_{t=1}^n R_{nt}^*}{\sqrt{\frac{1}{n} \sum_t R_{nt}^{*,2} + o_p(1)}} \\ &= \frac{\frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \sum_{t=j(a+b)+1}^{j(a+b)+a} R_{nt}^* + o_p(1)}{\sqrt{\frac{1}{n} \sum_{j=0}^{k-1} \left( \sum_{t=j(a+b)+1}^{j(a+b)+a} R_{nt}^* \right)^2 + o_p(1)}} = \frac{\frac{1}{\sqrt{n}} \sum_{j=0}^{k-1} \eta_j + o_p(1)}{\sqrt{\frac{1}{n} \sum_{j=0}^{k-1} \eta_j^2 + o_p(1)}} \Rightarrow N(0, 1) \end{aligned} \quad (77)$$

by Theorem 4.1 of de la Peña et al. (2009), since the denominator converges in probability to a strictly positive quantity  $(\sigma^2(x)\xi_2/\xi_1^{-2})$ , and that  $\eta_j$  belongs to the domain of attraction of a normal distribution by definition and (76). ■

### 7.3 Proof of Lemma 1 and 2

Lemma 1 is a straightforward extension of Lemma 4.3 and 4.4 of Ferraty and Vieu (2006), and hence is omitted. Lemma 2 can be shown by noting that for each  $n$  the  $\tau_n$ -dimensional polyhedron  $D := \{w = (w_i)_{i \leq \tau} \in \mathbb{R}^\tau, |w_i| \leq \lambda\}$  can be covered by  $(\lceil 2\lambda\sqrt{\tau}/\varepsilon \rceil + 1)^\tau$  number of balls of radius  $\varepsilon$ , see Chaté and Courbage (1997), and then following the proof steps of Theorem 2 in Jia et al. (2003). ■

### 7.4 Proof of Theorem 4

As before, we start from the decomposition (20):

$$\hat{m}(x) - m(x) = \frac{1}{\hat{m}_1(x)} \left( \left[ \hat{m}_2(x) - E\hat{m}_2(x) \right] + \left[ E\hat{m}_2(x) - m(x) \right] - m(x) \left[ \hat{m}_1(x) - 1 \right] \right).$$

We recall from (73) that  $\varphi_x(\lambda \underline{h}) \sim \varphi(\lambda \underline{h})$  and that the small deviation for the truncated regressor  $X = (X_1, \dots, X_\tau, 0, 0, \dots)$  denoted  $\varphi^T(\lambda \underline{h})$  satisfies

$$\varphi(\lambda \underline{h}) = P\left(\sum_{j=1}^{\infty} j^{-2p} X_j^2 \leq h^2\right) \leq P\left(\sum_{j=1}^{\tau} j^{-2p} X_j^2 \leq h^2\right) = \varphi^T(\lambda \underline{h}). \quad (78)$$

In the first step of the proof we show

$$\sup_{x \in \mathcal{S}_\tau} \left| \widehat{m}_2(x) - E\widehat{m}_2(x) \right| = O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda \underline{h})}}\right). \quad (79)$$

We cover the set  $\mathcal{S}_\tau$  defined in (41) with  $L = L(\mathcal{S}_\tau, \eta)$  number of balls of radius  $\eta$  denoted by  $I_k$ , each of which is centred at  $x_k$ ,  $k = 1, \dots, L$ . i.e.  $\mathcal{S}_\tau \subset \bigcup_{k=1}^{L_n} B(x_k, \eta)$ . Then it follows that

$$\begin{aligned} \sup_{x \in \mathcal{S}_\tau} \left| \widehat{m}_2(x) - E\widehat{m}_2(x) \right| &= \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \left| \widehat{m}_2(x) - E\widehat{m}_2(x) \right| \\ &= \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \left| \widehat{m}_2(x) - \widehat{m}_2(x_k) + \widehat{m}_2(x_k) - E\widehat{m}_2(x_k) + E\widehat{m}_2(x_k) - E\widehat{m}_2(x) \right| \\ &\leq \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \left| \widehat{m}_2(x) - \widehat{m}_2(x_k) \right| + \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \left| E\widehat{m}_2(x_k) - E\widehat{m}_2(x) \right| \\ &\quad + \max_{1 \leq k \leq L_n} \left| \widehat{m}_2(x_k) - E\widehat{m}_2(x_k) \right| =: R_1 + R_2 + R_3, \end{aligned} \quad (80)$$

where  $\widehat{m}_2(x_k) = (EK_1)^{-1} \sum_{t=1}^n Y_t K_{t,k}$  and  $K_{t,k} = K(\|H^{-1}(x_k - X_t)\|)$ .

We first consider  $R_1$ :

$$\begin{aligned} R_1 &= \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \left| \widehat{m}_2(x) - \widehat{m}_2(x_k) \right| \\ &= \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \left| \frac{1}{nEK_1} \sum_{t=1}^n Y_t K\left(\|H^{-1}(x - X_t)\|\right) - Y_t K\left(\|H^{-1}(x_k - X_t)\|\right) \right| \\ &\leq \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap \mathcal{S}_\tau} \frac{C}{n\varphi(\lambda \underline{h})} \sum_{t=1}^n |Y_t K_t - Y_t K_{t,k}|. \end{aligned}$$

Now, because the kernel function is assumed to be Lipschitz continuous by Assumption E2, it follows that

$$R_1 \leq \frac{1}{n} \sum_{t=1}^n \frac{C'|Y_t|}{\varphi(\underline{h}\lambda)} \eta h^{-1} =: \frac{1}{n\varphi(\underline{h}\lambda)} \sum_{t=1}^n J_t,$$

where  $J_t$  is  $\alpha$ -mixing under both assumptions A1' and A2'. Then for some  $\delta > 0$ , on choosing  $\eta = \log n/n$  and by Assumption B4' we see that the tail condition (which is required for the exponential inequality to be applied below) continues to hold for  $J_t$ .

Also, using Assumption B6 we see that

$$E|J_t| \leq \frac{E(E(|Y_t||X))\eta}{h} \leq \frac{C\eta}{h}. \quad (81)$$

By Lemma 2 we can specify the Kolmogorov's entropy for  $S_\tau$  with  $\eta = \log n/n^2$ :

$$\log L\left(S, \frac{\log n}{n^2}\right) = C \log \left[ \left( \frac{2\lambda n^2}{\sqrt{\log n}} + 1 \right)^{\log n} \right] \sim \log n \times \log \left[ \frac{2\lambda n^2}{\sqrt{\log n}} \right]$$

for sufficiently large  $n$  and  $\lambda$ , implying that the order of Kolmogorov's  $\frac{\log n}{n^2}$  entropy is

$$O\left(\log L\left(S_\tau, \frac{\log n}{n^2}\right)\right) = O\left((\log n)^2 - \log n[\log \log n]\right) = O\left((\log n)^2\right). \quad (82)$$

Now, we apply the Fuk-Nagaev inequality (cf. Fuk and Nagaev (1971)) for exponentially mixing variables of Merlevède, Peligrad and Rio (2009, 1.7) with  $\varepsilon = \varepsilon_0[\log L(S, \frac{\log n}{n^2})/(n\varphi(\lambda h))]^{1/2}$  and  $r = (\log n)^2$  for some positive constant  $\varepsilon_0$ . Since

$$s_n^2 := \sum_{t=1}^n \sum_{s=1}^n \text{cov}(J_t, J_s) \leq C \left( \frac{(\log n)^2}{n^2 h^2} \right) = O(n\varphi(\lambda h) \log n),$$

and due to exponential mixing assumption it follows that

$$\begin{aligned} & P\left( \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap S_\tau} |\hat{m}_2(x) - \hat{m}_2(x_k)| > \varepsilon_0 \sqrt{\frac{\log L(S, \frac{\log n}{n^2})}{n\varphi(\lambda h)}} \right) \\ & \leq 4 \left( 1 + \frac{n^2 \varphi(\lambda h)^2 \varepsilon_0^2 \log L(S, \frac{\log n}{n^2})}{16(\log n)^2 s_n^2 n\varphi(\lambda h)} \right)^{-\frac{(\log n)^2}{2}} + \frac{16Cn}{\sqrt{n\varphi(\lambda h)} \log n} e^{-\varsigma \left\{ \frac{\sqrt{n\varphi(\lambda h)}}{(\log n)} \right\}^\gamma} \\ & = 4 \left( 1 + \frac{n\varphi(\lambda h) \varepsilon_0^2}{16s_n^2} \right)^{-\frac{(\log n)^2}{2}} + \frac{16C\sqrt{n}}{\sqrt{\varphi(\lambda h)} \log n} \exp(-\varsigma \log L_n) \\ & \leq 4 \exp\left(-\frac{\varepsilon_0^2 (\log n)^2 n\varphi(\lambda h)}{32s_n^2}\right) + \left(\frac{Cn^2}{\sqrt{\log n}}\right) L_n^{-\varsigma} \\ & \leq 4 \exp\left(-\frac{\varepsilon_0^2 \log n}{32}\right) + \left(\frac{Cn^2}{\sqrt{\log n}}\right) \left(\frac{\sqrt{\log n}}{n^2}\right)^{\varsigma \log n} \longrightarrow 0 \end{aligned} \quad (83)$$

by the Taylor expansion of  $\log(1 + \epsilon)$  for sufficiently small  $\epsilon > 0$ .

Hence by (81) and Assumption E1 it follows that

$$\begin{aligned} R_1 &= \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap S_\tau} |\hat{m}_2(x) - \hat{m}_2(x_k)| \leq O\left(\frac{\eta}{h}\right) + O_P\left(\sqrt{\frac{\log L(S, \frac{\log n}{n^2})}{n\varphi(\lambda h)}}\right) \\ &= O\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda h)}}\right) + O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda h)}}\right) = O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda h)}}\right). \end{aligned} \quad (84)$$

As for the second term  $R_2$ , we have

$$R_2 \leq \max_{1 \leq k \leq L_n} \sup_{x \in I_k \cap S_\tau} E|\hat{m}_2(x) - \hat{m}_2(x_k)| = O\left(\frac{\eta}{h}\right) = O\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda h)}}\right). \quad (85)$$

Next we move on to the last component:

$$R_3 = \max_{1 \leq k \leq L_n} |\hat{m}_2(x_k) - E\hat{m}_2(x_k)| =: \max_{1 \leq k \leq L_n} |W_n(x_k)| \quad (86)$$

where

$$\begin{aligned} W_n(x) &= \hat{m}_2(x) - E\hat{m}_2(x) = \frac{1}{nEK_1} \sum_{t=1}^n [Y_t K_t - EY_t K_t] \\ &\leq \frac{C}{n\varphi_x^T(\underline{h}\lambda)} \sum_{t=1}^n [Y_t K_t - EY_t K_t] = \frac{C}{n\varphi_x(\underline{h}\lambda)} \sum_{t=1}^n U_{nt}. \end{aligned}$$

where  $U_{nt} = Y_t K_t - EY_t K_t$ , and by elementary arguments

$$P\left(\max_{1 \leq k \leq L_n} |\hat{m}_2(x_k) - E\hat{m}_2(x_k)| > \varepsilon\right) \leq L_n \cdot \sup_{x \in \mathcal{S}} P(|W_n(x)| > \varepsilon). \quad (87)$$

Due to the dependence of  $Q_{nt}$  on  $X_t$  we consider the cases of static and dynamic regressions separately, because the asymptotic arguments to follow depends upon the temporal dependence structure of  $Q_{nt}$ .

In the static case, we first examine the situation where the response is unbounded and satisfies the exponential tail condition in B4'. Since

$$s_n^2 = \sum_{t=1}^n \sum_{s=1}^n \text{cov}(U_{nt}, U_{ns}) = O(n\varphi_x^T(\underline{h}\lambda)),$$

we apply the Fuk-Nagaev inequality for exponentially mixing variables once again. Writing  $L_n := L(S, \frac{\log n}{n^2})$  and taking  $\varepsilon = \varepsilon_0 [\log L(S, \frac{\log n}{n^2}) / (n\varphi(\lambda \underline{h}))]^{1/2}$  and  $r = (\log n)^{2+\epsilon}$ ,  $\epsilon \in (0, 1/2)$  for some  $\varepsilon_0 > 0$ , we have

$$\begin{aligned} P\left(|\hat{m}_2(x) - E\hat{m}_2(x)| > \varepsilon_0 \sqrt{\frac{\log L_n}{n\varphi(\lambda \underline{h})}}\right) &\leq P\left(\left|\sum_{t=1}^n U_{nt}\right| > n\varphi_x^T(\underline{h}\lambda) \varepsilon_0 \sqrt{\frac{\log L_n}{n\varphi^T(\lambda \underline{h})}}\right) \\ &\leq 4 \left(1 + \frac{n\varphi^T(\lambda \underline{h}) \varepsilon_0^2 \log L_n}{16(\log n)^{2+\epsilon} s_n^2}\right)^{-\frac{(\log n)^{2+\epsilon}}{2}} + \frac{16Cn}{\sqrt{n\varphi^T(\lambda \underline{h})} \log n} \exp\left(-\varsigma \left\{\frac{\sqrt{n\varphi^T(\lambda \underline{h})}}{(\log n)^{1+\epsilon}}\right\}^\gamma\right) \\ &\leq 4 \left(1 + \frac{\varepsilon_0^2 \log L_n}{16(\log n)^{2+\epsilon}}\right)^{-\frac{(\log n)^{2+\epsilon}}{2}} + \frac{16C\sqrt{n}}{\sqrt{\varphi^T(\lambda \underline{h})} \log n} \exp(-\varsigma \log L_n) \\ &\leq 4 \exp\left(-\frac{\varepsilon_0^2 \log L_n}{32}\right) + \left(\frac{Cn^2}{\sqrt{\log n}}\right) L_n^{-\varsigma} \leq 4L_n^{-\frac{\varepsilon_0^2}{32}} + \left(\frac{Cn^2}{\sqrt{\log n}}\right) \left(\frac{\sqrt{\log n}}{n^2}\right)^{\varsigma \log n} \rightarrow 0 \end{aligned}$$

because  $\gamma \geq 1$  and  $L_n = O((n^2/\sqrt{\log n})^{\log n})$ .

Now since  $\varsigma > 1$ , by choosing  $\varepsilon_0$  large enough it follows by (87) that

$$R_3 = \max_{1 \leq k \leq L_n} |\hat{m}_2(x_k) - E\hat{m}_2(x_k)| = O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda \underline{h})}}\right). \quad (88)$$

In the special case when the response is bounded, the same result continues to hold with  $\gamma_1 = \infty$  (so that  $\gamma_2 = \gamma(\geq 1)$ ).

An alternative proof for the case of bounded response could be done by applying the exponential inequality of Bosq (1996, Theorem 1.3.2) for  $\alpha$ -mixing sequences as follows: Noting that  $|Q_t| \leq C/\varphi_x(\underline{h}\lambda) =: b$ ,  $\forall t$ , and that  $\sigma^2(r) := p \cdot \text{var}(Q_t) = O(p/\varphi(\underline{h}\lambda))$  (where  $p = n/(2q)$  and  $q = \log n \sqrt{n}/\sqrt{\varphi}$ ) by the Cauchy-Schwarz inequality and Assumption B4 we have

$$v^2(r) = \frac{2}{p^2} \sigma^2(r) + \frac{b\varepsilon}{2} \leq \frac{Cq}{n\varphi_x(\underline{h}\lambda)} + \frac{C\varepsilon}{\varphi_x(\underline{h}\lambda)} \leq \frac{C'\varepsilon}{\varphi_x(\underline{h}\lambda)},$$

where  $\varepsilon = \varepsilon_0 \sqrt{\log L_n / (n\varphi)}$  and  $L_n := L(S, \frac{\log n}{n^2})$ , and by Assumption A1 that

$$\begin{aligned} P\left(\left|\widehat{m}_2(x) - E\widehat{m}_2(x)\right| > \varepsilon_0 \sqrt{\frac{\log L_n}{n\varphi(\lambda\underline{h})}}\right) &\leq 4e^{-\varepsilon^2 q / (8v^2(r))} + 22\sqrt{1 + \frac{4b}{\varepsilon}} q \alpha\left(\left[\frac{n}{2q}\right]\right) \\ &\leq 4 \exp\left\{-\frac{\varepsilon_0 q \varphi \sqrt{\log L_n}}{8\sqrt{n\varphi}}\right\} + 22\left(1 + \frac{4\sqrt{n\varphi}}{\varphi \log n}\right)^{1/2} \frac{\log n \sqrt{n}}{\sqrt{\varphi}} \alpha\left(\left[\frac{\sqrt{n\varphi}}{2\log n}\right]\right) \\ &\leq 4 \exp\left\{-\frac{\varepsilon_0 \log L_n}{8}\right\} + \exp\left(-\varsigma \left\{\frac{\sqrt{n\varphi(\lambda\underline{h})}}{\log n}\right\}^{\gamma_2}\right) \\ &\leq 4L_n^{-\varepsilon_0/8} + \frac{C(\log n)^{1/2} n^{3/4}}{\varphi(\lambda\underline{h}) L_n^\varsigma} \rightarrow 0. \end{aligned}$$

In the dynamic regression case (i.e. under C2), the same conclusion can be derived by starting from (86) and exploiting the mixing approximation argument:

$$\begin{aligned} \max_{1 \leq k \leq L_n} |\widehat{m}_2(x_k) - E\widehat{m}_2(x_k)| &= \max_{1 \leq k \leq L_n} |W_n(x_k)| = \max_{1 \leq k \leq L_n} \left| n^{-1} \sum_{t=1}^n Q_{nt,k} \right| \\ &\leq \max_{1 \leq k \leq L_n} \left| \frac{1}{n} \sum_{t=1}^n Q_{nt,k,(r)} \right| + \sup_{x \in \mathcal{S}} \frac{1}{n} \sum_{t=1}^n |Q_{nt,(r)} - Q_{nt}| \\ &= O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda\underline{h})}}\right) + O_P\left(\frac{\sqrt{v_2(r)}}{\varphi(\lambda\underline{h})}\right) = O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda\underline{h})}}\right), \end{aligned}$$

since  $\sqrt{n} \sqrt{v_2(r)} (\log n)^{-1} / \sqrt{\varphi} \leq \sqrt{n} \sqrt{v_2(r)} / \varphi \rightarrow 0$  by (26).

Now returning back to where we started, viewing  $\widehat{m}_1(x)$  as a special case of  $\widehat{m}_2(x)$  with  $Y_t = 1$   $\forall t$ , we can repeat the above procedure, yielding (since  $E\widehat{m}_1(x) = 1$ )

$$\sup_{x \in \mathcal{S}_r} |\widehat{m}_1(x) - 1| = O_P\left(\sqrt{\frac{(\log n)^2}{n\varphi(\lambda\underline{h})}}\right). \quad (89)$$

The proof is now complete in view of (78), (79), (84), (85), (88), (89), contributions from the bias component, Proposition 4.1 of Dunker, Lifshits and Linde (1998), and Theorem 1.1 of Hong, Lifshits and Nazarov (2016). ■

## 7.5 Proof of Lemma 3

The main idea of the proof is to follow the arguments of Tsybakov (2004) where a lower bound for minimax risk of an estimator is constructed via deriving the upper bound of the squared Hellinger distance between probability measures. Further details are omitted as they closely follow the proof of Theorem 3 of Mas (2012). ■

## 7.6 Proof of Theorem 5

Given the extended moment condition upto  $8 + \delta$ , it is straightforward to see (from Theorem 1 and 2 & 3) the consistency of  $\hat{\sigma}^j(x_i)$  for  $\sigma^j(x_i)$  for  $j = 1, 2, 3, 4$  at every point of continuity  $x_i$ , and the asymptotic normality of  $(\hat{\mu}, \hat{\sigma}^2)$  with limiting variance  $\Omega(x_i)$ .

Hence it suffices to show asymptotic independence of  $\hat{m}(x_i)$  and  $\hat{m}(x'_i)$  across  $i$ , where  $x_i$  and  $x'_i$  are continuity points of  $m$  such that  $\|D^{-1}(x_i - x'_i)\| > 0$ . Following the notations of the proof of Theorem 2 and 3, the asymptotic covariance matrix is given by  $\text{Var}[(\sqrt{\phi(h)}/\sqrt{n}) \sum_{t=1}^n R_{nt}]$ , and

$$\text{Var}(R_{nt}) = \text{Var} \begin{pmatrix} \frac{1}{EK_{1,x}} \cdot K_{t,x}[Y_t - m(x)] \\ \frac{1}{EK_{1,x'}} \cdot K_{t,x'}[Y_t - m(x')] \end{pmatrix} = E \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \quad (90)$$

We know from Theorem 2 and 3 that as for  $A_{11} \simeq \sigma^2(x)$  and  $A_{22} \simeq \sigma^2(x')$ . So we just consider the off-diagonal terms. Due to stationarity we see that

$$\begin{aligned} & E[K_{t,x}K_{t,x'}(Y_t - m(x))(Y_t - m(x'))] \\ &= E[K_{1,x}K_{1,x'} \left\{ (Y_1 - m(X_1)) + (m(X_1) - m(x)) \right\} \left\{ (Y_1 - m(X_1)) + (m(X_1) - m(x')) \right\}] \\ &= E[K_{1,x}K_{1,x'}(Y_1 - m(X_1))(Y_1 - m(X_1))] + o(1) = E[K_{1,x}K_{1,x'}\sigma^2(X_1)] + o(1) \\ &\leq \sup_{u \in B(x,h) \cap B(x',h)} \sigma^2(u) E[K_{1,x'}K_{1,x}] \rightarrow 0 \end{aligned}$$

as  $h \rightarrow 0$  since the kernels return 0 outside its compact support and  $\|D^{-1}(x_i - x'_i)\| > 0$ . The desired result now directly follows via the delta method. ■

## References

- [1] Andrews, D. W. K. (1984). *Non-Strong Mixing Autoregressive Processes*. Journal of Applied Probability, 21(4), 930-934.
- [2] Andrews, D. W. K. (1995). *Nonparametric kernel estimation for semiparametric models*. Econometric Theory, 11(3), 560-586.
- [3] Azais, J. M. and Fort, J. C. (2013). *Remark on the finite-dimensional character of certain results of functional statistics*. Comptes Rendus Mathematique, 351(3), 139-141.

- [4] Backus, K. and Gregory, A. W. (1993). *Theoretical Relations Between Risk Premiums and Conditional variances*. Journal of Business and Economic Statistics, 11(2), 177-185.
- [5] Bansal, R. and Yaron, A. (2004): *Risks for the Long Run: A Potential Resolution of Asset Pricing Puzzles*. Journal of Finance, 59(4), 1481-1509.
- [6] Bierens, H. J. (1983). *Uniform consistency of kernel estimators of a regression function under generalized conditions*. Journal of the American Statistical Association, 78(383), 699-707.
- [7] Bierens, H. J. (1987). *Kernel estimators of regression functions*. In Advances in econometrics: Fifth world congress, 99-144, Vol. 1, (ed.) Bewley, T. F., Cambridge: Cambridge University Press.
- [8] Bingham, N. H., Goldie, C. M. and Teugels, J. L. (1987). *Regular variation*. Cambridge: Cambridge University Press.
- [9] Billingsley, P. (1968). *Convergence of Probability Measures*. New York: John Wiley.
- [10] Borovkov, A. A. and Ruzankin, P. S. (2008). *On small deviations of series of weighted random variables*. Journal of Theoretical Probability, 21(3), 628-649.
- [11] Bosq, B. (1996). *Nonparametric statistics for stochastic processes: estimation and prediction*. New York: Springer-Verlag.
- [12] Bradley, R. C. (2005). *Basic properties of strong mixing conditions. A survey and some open questions*. Probability Surveys, 2(2), 107-144.
- [13] Campbell, J. Y. and Hentschel, L. (1992). *No News is Good News: An Asymmetric Model of Changing Volatility in Stock Returns*. Journal of Financial Economics, 31(3), 281-318.
- [14] Campbell, J. Y. (1987). *Stock Returns and the Term Structure*. Journal of Financial Economics, 18(2), 373-399.
- [15] Campbell, J. Y. and Cochrane, J. H. (1999). *By Force of Habit: A Consumption-Based Explanation of Aggregate Stock Market Behavior*. Journal of Political Economy, 107(2), 205-251.
- [16] Chaté, H. and Courbage, M. (1997). *Special issue on lattice dynamics*. Physica D, 103, 1-611.
- [17] Chen, R. and Tsay, R. S. (1993). *Nonlinear additive ARX models*. Journal of the American Statistical Association, 88(423), 955-967.
- [18] Chen, X., and T.M. Christensen (2015) Optimal uniform convergence rates and asymptotic normality for series estimators under weak dependence and weak conditions. Journal of Econometrics 188, 447-465.

- [19] Conrad, C. and Mammen, E. (2008), *Nonparametric regression on latent covariates with an application to semiparametric GARCH-in-Mean models*. Discussion Paper No. 473, University of Heidelberg, Department of Economics.
- [20] Constantinides, G. (1990). *Habit Formation: A Resolution of the Equity Premium Puzzle*. Journal of Political Economy, 98(3), 519-543.
- [21] Cox, J., Ingersoll, J. E. and Ross, S. A. (1985). *An Intertemporal General Equilibrium Model of Asset Prices*. Econometrica, 53(2), 363-384.
- [22] Davidson, J. (1994). *Stochastic Limit Theory: An Introduction for Econometricians*. Oxford: Oxford University Press.
- [23] Davydov, Y. A. (1968). *Convergence of distributions generated by stationary stochastic processes*. Theory of Probability and Its Applications, 13(4), 691-696.
- [24] Delsol, L. (2009). *Advances on asymptotic normality in non-parametric functional time series analysis*. Statistics, 43(1), 13-33.
- [25] Devroye, L. P. (1978). *The uniform convergence of the Nadaraya-Watson regression function estimate*. Canadian Journal of Statistics, 6(2), 179-191.
- [26] Devroye, L. (1981). *On the almost everywhere convergence of nonparametric regression function estimates*. Annals of Statistics, 9(6), 1310-1319.
- [27] Doukhan, P. (1994). *Mixing*. New York: Springer.
- [28] Doukhan, P. and Wintenberger, O. (2008). *Weakly dependent chains with infinite memory*. Stochastic Processes and their Applications, 118(11), 1997-2013.
- [29] Duflo, M. (1997). *Random iterative models*. Berlin: Springer-Verlag.
- [30] Dunker, T., Lifshits, M. A. and Linde, W. (1998). *Small deviation probabilities of sums of independent random variables*. In: Eberlein, E. (ed.), High dimensional probability, volume 43 of Progress in Probability., Birkhauser, Basel, 59-74.
- [31] Engle, R. F., Lilien, D. M. and Robins, R. P. (1987). *Estimating Time varying Risk Premia in the Term Structure: The ARCH-M Model*. Econometrica, 55(2), 391-407.
- [32] Escanciano, J.C., Pardo-Fernández, J.C. and Van Keilegom, I. (2015). *Semiparametric estimation of risk return relationships*. Journal of Business and Economic Statistics (To Appear).
- [33] Fan, J. and Masry, E. (1992). *Multivariate regression estimation with errors-in-variables: asymptotic normality for mixing processes*. Journal of Multivariate Analysis, 43(2), 237-271.



- [34] Fan, J. (1990). *A remedy to regression estimators and nonparametric minimax efficiency*. Technical Report 161, Department of Statistics, University of North Carolina at Chapel Hill.
- [35] Fan, J. and Yao, Q. (2005). *Nonlinear Time Series: Nonparametric and Parametric Methods*. New York: Springer.
- [36] Feller, W. (1971). *Introduction to Probability Theory and Its Applications*, Vol. 2. New York: Wiley.
- [37] Ferraty, F., Laksaci, A., Tadj, A. and Vieu, P. (2010). *Rate of uniform consistency for nonparametric estimates with functional variables*. Journal of Statistical Planning and Inference, 140(2), 335-352.
- [38] Ferraty, F. and Romain, Y. (2010). *The Oxford Handbook of Functional Data Analysis*. New York: Oxford University Press.
- [39] Ferraty, F. and Vieu, P. (2002). *The functional nonparametric model and application to spectrometric data*. Computational Statistics, 17(4), 545-564.
- [40] Ferraty, F. and Vieu, P. (2006). *Nonparametric functional data analysis: Theory and Practice*. New York: Springer.
- [41] French, K. R., Schwert, G. W. and Stambaugh, R. F. (1987) *Expected Stock Returns and Volatility*. Journal of Financial Economics, 19(1), 3-29.
- [42] Fuk, D. K. and Nagaev, S. V. (1971). *Probability inequalities for sums of independent random variables*. Theory of Probability and its applications, 16(4), 643-660.
- [43] Gao, F., Hannig, J. and Torcaso, F. (2003). *Comparison Theorems for Small Deviations of Random Series*. Electronic Journal of Probability, 8(21), 1-17.
- [44] Gallant, A. R., Hsieh, D. and Tauchen, G. (1989). *On Fitting a Recalcitrant Series: The Dollar/Pound Exchange Rate, 1974-1983*. In W. A. Barnett, J. Powell, and G. Tauchen (eds.), *Nonparametric and Semiparametric Methods in Econometrics and Statistics*.
- [45] Gennotte, G. and Marsh, T. A. (1993). *Variations in economic uncertainty and risk premiums on capital assets*. European Economic Review, 37(5), 1021-1041.
- [46] Geenens, G. (2011). *Curse of dimensionality and related issues in nonparametric functional regression*. Statistics Surveys, 5, 30-43.
- [47] Ghysels, E., Santa-Clara, P. and Valkanov, R. (2005). *There is a Risk-Return Trade-Off After All*. Journal of Financial Economics, 76(3), 509-548.

- [48] Glosten, L., Jagannathan, R. and Runkle, D. E. (1993). *On The Relation Between the Expected Value and the Volatility of the Nominal Excess Return on Stocks*. Journal of Finance, 48(5), 1779-1801.
- [49] Götze, F. and Hipp, C. (1994). *Asymptotic distribution of statistics in time series*. Annals of Statistics, 22(4), 2062-2088.
- [50] Greblicki, W. and Krzyzak, A. (1980). *Asymptotic properties of kernel estimates of a regression function*. Journal of Statistical Planning and Inference, 4(1), 81-90.
- [51] Guo, H. and Whitelaw, R. (2006). *Uncovering the Risk-Return Relation in the Stock Market*. Journal of Finance, 61(3), 1433-1463.
- [52] Härdle, W. K. (1990). *Applied Nonparametric Regression*. Cambridge: Cambridge University Press.
- [53] Harvey, C. (2001). *The Specification of Conditional Expectations*. Journal of Empirical Finance, 8(5), 573-638.
- [54] Hong, S. Y., Lifshits, M. and Nazarov. A. (2016). *Small deviations in  $L_2$ -norm for Gaussian dependent sequences*. Electronic Communications in Probability, 21(41), 1-9.
- [55] Ibragimov, I. A. (1962). *Some limit theorems for stationary processes*. Theory of Probability & Its Applications, 23, 291-300.
- [56] Ibragimov, I. A. and Linnik, Y. B. (1971). *Independent and Stationary Sequences of Random variables*. Groningen: Wolters-Noordhoff.
- [57] Jia, Q., Zhou, S. and Yin. F. (2003). *Kolmogorov entropy of global attractor for dissipative lattice dynamical systems*. Journal of Mathematical Physics, 44, 5804-5801.
- [58] Karamata, J. (1933). *Sur un mode de croissance regulire. Theoremes fondamentaux* (in French). Bulletin de la Societe Mathematique de France, 61, 55-62.
- [59] Lettau, M. and Ludvigson, S. (2003). *Measuring and modeling variation in the risk-return tradeoff*. In: Ait-Sahalia, Y., Hansen, L. (Eds.), Handbook of Financial Econometrics. Amsterdam: North-Holland.
- [60] Lettau, M. and Nieuwerburgh, S. V. (2008). *Reconciling the Return Predictability Evidence*. Review of Financial Studies, 21(4), 1607-1652.
- [61] Li, W. V., 2012. *Small value probabilities in analysis and mathematical physics*. Presented at the Arizona School of Analysis and Mathematical Physics, Tucson, Arizona, United States in March 15, 2012. Retrieved from the link [http://math.arizona.edu/~mathphys/school\\_2012/WenboLi.pdf](http://math.arizona.edu/~mathphys/school_2012/WenboLi.pdf).

- [62] Li, W. V. and Shao, Q. M. (2001). *Gaussian processes: inequalities, small ball probabilities and applications*. Handbook of Statistics, 19, 533-598.
- [63] Linton, O. and Perron, B. (2003). *The Shape of the Risk Premium: Evidence From a Semiparametric Generalized Autoregressive Conditional Heteroscedasticity Model*. Journal of Business and Economic Statistics, 21(3). 354-367.
- [64] Linton, O. B. and Sancetta, A. (2009). *Consistent estimation of a general nonparametric regression function in time series*. Journal of Econometrics, 152(1), 70-78.
- [65] Lu, Z. (2001). *Asymptotic normality of kernel density estimators under dependence*. Annals of the Institute of Statistical Mathematics, 53(3), 447-468.
- [66] Ludvigson, S. and Ng, S. (2007) *The Empirical Risk-Return Relation: A Factor Analysis Approach*. Journal of Financial Economics, 83(1), 171-222.
- [67] Lundblad, C. (2005). *The Risk-Return Tradeoff in the Long Run: 1836-2003*. Journal of Financial Economics, 85(1), 123-150.
- [68] Mas, A. (2012). *Lower bound in regression for functional data by small ball probability representation in Hilbert space*, Electronic Journal of Statistics, 6, 1745-1778.
- [69] Masry, E. (2005). *Nonparametric regression estimation for dependent functional data: asymptotic normality*. Stochastic Processes and their Applications, 115(1), 155-177.
- [70] Masry, E. and Fan, J. (1997). *Local polynomial estimation of regression function for mixing processes*. Scandinavian Journal of Statistics, 24(2), 1965-1979.
- [71] Merlevède, F., Peligrad, M., Rio, E., 2009. *A Bernstein type inequality and moderate deviations for weakly dependent sequences*. Preprint available at <http://arxiv.org/abs/0902.0582>.
- [72] Merton, R. C. (1973). *An Intertemporal Capital Asset Pricing Model*. Econometrica, 41(5), 867-887.
- [73] Nadaraya, E. A. (1964). *On estimating regression*. Theory of Probability & Its Applications, 9(1), 141-142.
- [74] Nadaraya, E. A. (1970). *Remarks on non-parametric estimates for density functions and regression curves*. Theory of Probability & Its Applications, 15(1), 134-137.
- [75] Olver, F. W. J., Lozier, D. W., Boisvert, R. F. and Clark, C. W. (2010). *NIST Handbook of Mathematical Functions*. New York: Cambridge University Press.

- [76] Pagan, A. R., and Hong, Y. S. (1990). *Nonparametric Estimation and the Risk Premium*. In W. A. Barnett, J. Powell, and G. Tauchen (eds.), *Nonparametric and Semiparametric Methods in Econometrics and Statistics*.
- [77] Pagan, A. and Ullah, A. (1988). *The econometric analysis of models with risk terms*. Journal of Applied Econometrics, 3(2), 87-105.
- [78] Parzen, E. (1962). *On estimation of a probability density function and mode*. Annals of Mathematical Statistics, 33(3), 1065-1076.
- [79] Pástor, L., Sinha, M. and Swaminathan, B. (2008). *Estimating the Intertemporal Risk-Return Tradeoff Using the Implied Cost of Capital*. Journal of Finance, 63(6), 2859-2897.
- [80] de la Peña, V. H., Lai, T. L. and Shao, Q. M. (2009). *Self-normalized processes: Limit theory and Statistical Applications*. New York: Springer.
- [81] Phillips, P. C. and Park, J. Y. (1998). *Nonstationary density estimation and kernel autoregression*. Cowles Foundation discussion paper.
- [82] Ramsay, J. O. and Silverman, B. W. (2002). *Applied functional data analysis: methods and case studies*. New York: Springer.
- [83] Rio, E. (2000). *Theorie asymptotique des processus aléatoires faiblement dépendants* (in French). Berlin: Springer Verlag.
- [84] Robinson, P. M. (1983). *Nonparametric estimators for time series*. Journal of Time Series Analysis, 4(3), 185-207.
- [85] Rosenblatt, M. (1956). *A central limit theorem and a strong mixing condition*. Proceedings of the National Academy of Sciences, 42(1), 43-47.
- [86] Roussas, G. G. (1989). *Consistent regression estimation with fixed design points under dependence conditions*. Statistics & Probability Letters, 8(1), 41-50.
- [87] Roussas, G. G. (1990). *Nonparametric regression estimation under mixing conditions*. Stochastic Processes and Their Applications, 36(1), 107-116.
- [88] Schuster, E. F. (1972). *Joint asymptotic distribution of the estimated regression function at a finite number of distinct points*. Annals of Mathematical Statistics, 43(1), 84-88.
- [89] Scruggs, J. (1998). *Resolving the Puzzling Intertemporal Relation Between the Market Risk Premium and the Conditional Market variance: A Two-Factor Approach*. Journal of Finance, 53(2), 575-603.

- [90] Scruggs, J. and Glabadanidis, P. (2003). *Risk Premia and the Dynamic covariance Between Stock and Bond Returns*. Journal of Financial and Quantitative Analysis, 38(2), 295-316.
- [91] Skorohod, A. B. (1967). *On the densities of probability measures in functional spaces*. Proc. Fifth Berkeley Sympos. Math. Statist. and Probability (Berkeley, California, 1965/66) Vol. II: Contributions to Probability Theory, 1, 163-182. University of California Press, Berkely, CA.
- [92] Stone, C. (1980). *Optimal rates of convergence for nonparametric estimators*. Annals of Statistics, 8(6), 1348-1360.
- [93] Stone, C. (1982). *Optimal global rates of convergence for nonparametric regression*. Annals of Statistics, 10(4), 1040-1053.
- [94] Sytaya, G. N. (1974). *On certain asymptotic representations for a Gaussian measure in Hilbert space* (in Russian). Theory of Random Process, 2, 93-104.
- [95] Tsybakov, A. (2009). *Introduction to Nonparametric Estimation*. New York: Springer.
- [96] Watson, G. S. (1964). *Smooth regression analysis*. Sankhya Series A, 26(4), 359-372.
- [97] Whitelaw, R. (1994). *Time variations and covariations in the Expectation and Volatility of Stock Market Returns*, Journal of Finance, 49(2), 515-541.
- [98] Wu, W. B. (2011). *Asymptotic theory for stationary processes*. Statistics and Its Interface, 4, 207-226.
- [99] Zhang, X., M.L. King, and R.J. Hyndman (2006). A Bayesian approach to bandwidth selection for multivariate kernel density estimation. Computational Statistics and Data Analysis 50, 3009-3031.
- [100] Zolotarev, V. M. (1986). *Asymptotic behavior of the Gaussian measure in  $\ell_2$* . Journal of Soviet Mathematics, 35(2), 2330-2334.