



Semi-parametric estimation of the autoregressive parameter in non-Gaussian Ornstein–Uhlenbeck processes

S. Rao Jammalamadaka^a and Emanuele Taufer^b

^aDepartment of Statistics and Applied Probability, University of California, Santa Barbara, Santa Barbara, CA, USA;

^bDepartment of Economics and Management, University of Trento, Trento, Italy

ABSTRACT

This paper considers the problem of estimating the autoregressive parameter in discretely observed Ornstein–Uhlenbeck processes. Two consistent estimators are proposed: one obtained by maximizing a kernel-based likelihood function, and another by minimizing a Kolmogorov-type distance from independence. After establishing the consistency of these estimators, their finite-sample performance and possible normality in large samples, is investigated by means of extensive simulations. An illustrative example to credit rating is discussed.

ARTICLE HISTORY

Received 18 December 2017

Accepted 6 April 2018

KEYWORDS

Adaptive estimation; Kernel density estimation; Lévy process; self-decomposable distribution; minimum squared distance to independence

1. Introduction

A continuous stationary process $\{X(t), t \geq 0\}$ is defined to be of the Ornstein–Uhlenbeck type (OU for short) if it is the solution of the stochastic differential equation

$$dX(t) = -\lambda X(t)dt + d\dot{Z}(t) \quad (1)$$

here $\lambda > 0$, and $\dot{Z}(t)$ is a homogeneous Lévy process, commonly referred to as the background driving Lévy process (BDLP), which satisfies the condition $E[\log(1 + |\dot{Z}(1)|)] < \infty$ (see, e.g., Barndorff-Nielsen and Shephard 2001). Modeling via the use of general Lévy processes, other than Brownian motion, allows one to introduce specific non-Gaussian distributions for the marginal law of $X(t)$, and has received considerable attention in recent literature in an attempt to accommodate features such as jumps, semi-heavy tails and asymmetry, which are quite evident in real phenomena and are of practical interest in fields of application such as finance and econometrics.

Most notable examples include OU processes with marginal distributions such as the normal inverse Gaussian and the inverse Gaussian (Barndorff-Nielsen 1998), the variance gamma (Seneta 2004), the Meixner (Schoutens and Teugels 1998), the t-distribution (Heyde and Leonenko 2005), the normal, the stable and the gamma distributions. OU processes with positive jumps with marginal distributions such as the inverse Gaussian are often used as building blocks in stochastic volatility models (see, e.g., Barndorff-Nielsen and Shephard 2001).

A key concept related to these processes is that of self-decomposability. Recall that a random variable X with characteristic function $\psi(\zeta)$, is said to be self-decomposable if, for

CONTACT Emanuele Taufer emanuele.taufer@unitn.it Department of Economics and Management, University of Trento, Trento, Italy.

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/lssp.

© 2018 Taylor & Francis Group, LLC

47 all $c \in (0, 1)$, there exists a characteristic function $\psi_c(\zeta)$ such that $\psi(\zeta) = \psi(c\zeta)\psi_c(\zeta)$.
 48 Self-decomposability is closely related to stationary linear autoregressive time series of order
 49 1, i.e. an AR(1) process: essentially the only possible AR(1) processes are those for which
 50 the one-dimensional marginal law is self-decomposable and similarly for the OU process, i.e.
 51 an “AR(1)” in continuous time. For further details on self-decomposable, infinitely divisible
 52 distributions and Lévy processes see Sato (1999).

53 This paper is concerned with estimation of the autoregressive parameter λ . Maximum
 54 likelihood estimation of λ is generally infeasible except for a few special cases and the large
 55 availability of marginal distributions for X calls for efficient estimation in a broad range of
 56 situations. To this end we propose two estimators: an estimator using a kernel estimate of the
 57 likelihood; another based on minimum distance from independence, which addresses some
 58 of the problems that the kernel-based estimator encounters in certain cases.

59 Suppose we observe the process Eq. (1) at equi-spaced time points $0 < t_1 < \dots < \dots t_n$
 60 with $\Delta = t_j - t_{j-1}$, $j = 1, \dots, n$, $t_0 = 0$. In order to slightly simplify notation, denote the
 61 observation at time t_j , $X(t_j)$, by X_j . It follows from the discussion in Wolfe (1982) that, for self-
 62 decomposable distributions, a discrete AR(1) process can be embedded into a continuous OU
 63 process. In our case, this amounts to saying that the discretely observed OU process Eq. (1)
 64 can be written as

$$65 \quad X_j = e^{-\lambda\Delta}X_{j-1} + \varepsilon_j, \quad j = 1, 2, \dots, n \quad (2)$$

66 where the ε_j 's are *i.i.d.* random variables. Note that in practical applications, determining
 67 the timing of observations is quite arbitrary, which amounts to saying that from a practical
 68 point of view one is not able to distinguish between Δ and λ . In this paper, contrary to other
 69 approaches where Δ is assumed to be known, we will actually consider estimation of, say,
 70 $\lambda' = \lambda\Delta$ so that, from now on it will be assumed without loss of generality that $\Delta = 1$.
 71 Denote $\theta = e^{-\lambda\Delta}$ and rewrite Eq. (2) as

$$72 \quad X_j = \theta X_{j-1} + \varepsilon_j, \quad \theta \in \Theta, \quad \Theta = (0, 1), \quad j = 1, 2, \dots, n \quad (3)$$

73 With X_0 having distribution corresponding to the characteristic function $\psi(\zeta)$, model Eq. (3)
 74 is strictly stationary with marginal distribution having characteristic function $\psi(\zeta)$ and *i.i.d.*
 75 innovations with characteristic function $\psi_\theta(\zeta) = \psi(\zeta)/\psi(\theta\zeta)$.

76 Estimation of these models and in particular the estimation of the parameter θ (or λ) has
 77 attracted considerable interest in recent literature. When X is normal, the sample counterpart
 78 of the auto-correlation $Cor(X_1, X_2)$ provides, after transformation, the maximum likelihood
 79 estimator of λ . This turns out to be an estimator widely used in practice; Long (2009) has
 80 shown that the auto-correlation (AC) estimator is consistent for the model Eq. (3) with stable
 81 innovations with index of stability $1 < \alpha < 2$ with $\Delta = \Delta_n = 1/n$ when $n \rightarrow \infty$ and
 82 dispersion approaching 0. Zhang and Zhang (2013) show that the AC-based estimator of λ
 83 to be consistent for symmetric α -stable innovations for $0 < \alpha < 2$ either for fixed Δ and
 84 $\Delta \rightarrow 0$. Again, Hu and Long (2009) consider a least squares estimator for the case of α -
 85 stable innovations and show its consistency for $1 < \alpha < 2$ and $\Delta \rightarrow 0$. These approaches
 86 are equivalent when $\Delta \rightarrow 0$. Notwithstanding, the AC estimator turns out to be inefficient in
 87 many non-normal cases; to correct this situation Koul (1986) introduced a class of L_2 -distance
 88 estimators of θ when the errors have an unknown symmetric distribution.

93 Jongbloed, Van der Meulen, and Van der Waart (2005) have proposed a highly efficient
 94 estimator of θ for the case where model \dot{Z} is a subordinator, i.e., a process with positive
 95 increments. In this case, for the discretely observed model Eq. (3), $\hat{\theta} = \min_{1 \leq j \leq n} X_j/X_{j-1}$
 96 which had also been discussed by Nielsen and Shephard (2003) in a model with exponential
 97 innovations. For other estimation problems for non-negative Lévy-driven OU processes, see
 98 Brockwell, Davis, and Yang (2007).

99 Restricting attention to non-negative Lévy-driven OU processes, however, excludes a
 100 whole range of possible marginal distributions for the model Eq. (1). A general paramet-
 101 ric approach is considered by Taufer and Leonenko (2009a) which uses the characteristic
 102 function to estimate θ together with the parameters of the marginal distribution of X , while
 103 Andrews, Calder, and Davis (2009) discuss estimation of α -stable auto-regressive processes;
 104 see also Taufer, Leonenko, and Bee (2011) and Meintanis and Taufer (2012) for extensions
 105 to stochastic volatility models. Other papers of interest here are those of Diop and Yode
 106 (2010) who study a minimum distance estimator of θ when dispersion of the innovations
 107 approaches 0, and Ma (2010) who shows that the results of Long (2009) hold also under
 108 weaker conditions and Zhang, Lin, and Zhang (2015) which discuss LSE estimation for Lévy-
 109 driven moving averages.

110 The problem discussed here is closely connected to the works on adaptive estimation;
 111 in particular of direct relevance here are the papers of Kreiss (1987), Drost, Klaassen, and
 112 Werker (1997), Koul and Schick (1997) and Hallin et al. (2000) in time series contexts;
 113 Linton and Xiao (2007), Linton, Sperlich, and Van Keilegom (2008), and Yao and Zhao
 114 (2013) in semi-parametric and regression contexts; these approaches have in common the
 115 requirement that a preliminary consistent estimator of the parameter of interest is available
 116 while the approach proposed here has a one-step structure without using any preliminary
 117 estimator: only Eq. (3) is exploited and a simple maximization of a kernel density estimator is
 118 required. In this sense, the paper closest to our approach is Yuan and De Gooijer (2007) which
 119 considers a one-step adaptive procedure in the regression context. The problem discussed
 120 here may be seen as an extension to the dependent case of Yuan and De Gooijer (2007),
 121 although we adopt different techniques and require a minimal set of conditions, such as not
 122 requiring symmetry and placing very mild moment conditions in proving consistency of the
 123 estimators which generally hold in a large variety of self-decomposable distributions for OU
 124 processes.

125 As for our second estimator based on the minimum distance to independence, previous
 126 related literature dates back to Manski (1983), Brown and Wegkamp (2002), and Linton,
 127 Sperlich, and Van Keilegom (2008); in these papers a minimum mean squared distance to
 128 independence is considered; this approach would require the existence of a finite mean, while
 129 here with the aim of requiring a minimal set of conditions, we use a Kolmogorov-type distance
 130 instead; the use of this distance has been discussed by Manski (1983) for the case where a
 131 parametric form of the distribution function is given while here a semi-parametric setting is
 132 discussed.

133 In the next section we will precisely define the estimators and present the main results.
 134 In Section 3 the small sample performance of the estimators will be analyzed by means of
 135 extensive simulations. An Appendix presents the proofs of the results.

136 In the paper we will use the notation $X_n = O_p(a_n)$ meaning that, for any $\varepsilon > 0$ there exists
 137 a finite M such that $P(|X_n/a_n| > M) < \varepsilon \forall n$ and $X_n = o_p(a_n)$ meaning that, for any $\varepsilon > 0$,
 138 $\lim_{n \rightarrow \infty} P(|X_n/a_n| > \varepsilon) = 0$. $X_n \rightarrow_D X$ is used to indicate convergence in distribution.

139

2. Semi-parametric estimators for θ and main results

140

If we denote by $\theta_0 \in \Theta$ the true parameter value, then the sequence of innovations $\varepsilon_j = X_j - \theta_0 X_{j-1}$, $j = 1, 2, \dots, n$ is *i.i.d.* More generally, define the residuals $e_j = e_j^\theta = X_j - \theta X_{j-1}$, $j = 1, 2, \dots, n$. Note that only the choice $\theta = \theta_0$ assures that X is strictly stationary with *i.i.d.* innovations ε ; other choices of θ will lead to dependent innovations e . In fact, writing $e_j = (\theta_0 - \theta)X_{j-1} + \varepsilon_j$, $j = 1, 2, \dots, n$ we note that the sequence of the e_j 's is not independent due to the dependence of the X_j 's.

147

148

2.1. A Kernel-based estimator

149

Let $f_\theta = f_\theta(e)$ denote the density of the residuals and, for $\theta = \theta_0$, $f_{\theta_0} = f_{\theta_0}(\varepsilon)$ denotes the density of the innovations. Also, let $f_\theta(x_0, x_1)$ be the bivariate density of X_0 and X_1 .

151

Define the kernel estimator of f_θ , based on e_j , $j = 1, 2, \dots, n$, as

152

153

154

155

$$\hat{f}_\theta(x) := \frac{1}{nh} \sum_{j=1}^n K\left(\frac{x - e_j^\theta}{h}\right) \quad (4)$$

where K is a scalar kernel and $h = h(n)$ is a bandwidth sequence. The following estimator of θ is proposed:

157

158

159

160

$$\hat{\theta}_1 = \arg \max_{\theta \in \Theta} \sum_{i \in S} \log \hat{f}_\theta(e_i) := \arg \max_{\theta \in \Theta} L_n(\theta) \quad (5)$$

161

162

163

164

165

where S is a subset of $\{1, 2, \dots, n\}$ and it is introduced in case it is felt necessary to trim out some summands. Typically S will coincide with the full set $\{1, 2, \dots, n\}$, i.e. all observations are used to estimate f_θ however, in some instances, one could get very small positive estimates of \hat{f}_θ which can cause numerical problems due to un-boundedness of the logarithmic function near the origin. Also, negative estimates of \hat{f}_θ could arise if higher order kernels are used.

166

167

168

169

To avoid these problems it is quite common in entropy estimation to assume that the support of f_θ is bounded, see, e.g. Hall (1986), van Es (1992), Hall and Morton (1993), and Yuan and De Gooijer (2007). This is not the approach followed here where OU processes require unbounded distributions.

170

171

172

From a purely practical point of view it might be a sensible precaution to exclude those e_i such that $\hat{f}_\theta(e_i) < b$ for some prescribed positive b or, alternatively, omitting those e_i such that $|e_i| > M$.

173

174

175

176

177

178

In our simulations (Section 3) in all but the stable cases with index of stability less than 1 the whole set of data was used without noticing any problem. When some trimming is necessary, this is usually quite evident as one gets unreasonable estimates of θ , i.e. 0 or 1 or a seemingly unbounded numerical likelihood. In the simulations we have set some common level of trimming for a given distribution. For an actual application, close inspection of data and estimates would suggest which data values should be excluded from the computations.

179

180

181

From a theoretical point of view, in order to ensure consistency one needs to allow arbitrary values of b or M . This point will be discussed more fully in the appendix. For showing the consistency of $\hat{\theta}_1$, we need the following standard conditions in kernel estimation:

182

183

184

A1 The sequence $\{X_i\}_{0 \leq i \leq n}$ follows model Eq. (3) and is strictly stationary with non-degenerate self-decomposable marginal distribution such that, for some $p > 0$ $E(X_0^p) < \infty$.

- 185 A2 The density $f_\theta(x)$ is bounded away from 0 and Lipschitz continuous wrt θ on compact
 186 intervals of $x \in \mathbb{R}$ and $\sup_e f_\theta(e) < \infty$ for any θ .
 187 A3 The joint density $f_\theta(x_0, x_1)$, is bounded away from 0 on compact sets of $x_0, x_1 \in \mathbb{R}$ and
 188 $\sup_{x_0, x_1} f_\theta(x_0, x_1) < \infty$ for any θ .
 189 A4 $\int_{-\infty}^{\infty} |\log f_\theta(x)| f_\theta(x) dx < \infty$ for any θ .
 190 A5 $|K(u)| < \infty$, $\int_{-\infty}^{\infty} |K(u)| du < \infty$, $\int_{-\infty}^{\infty} |uK(u)| du < \infty$.
 191 A6 For some $M_1 < \infty$ and $M_2 < \infty$, either $K(u) = 0$ for $|u| > M_2$ and for all $u, u' \in \mathbb{R}$,
 192 $|K(u) - K(u')| \leq M_1|u - u'|$ or $K(u)$ is differentiable, $|(\partial/\partial u)K(u)| \leq M_1$, and for
 193 some $\nu > 1$, $|(\partial/\partial u)K(u)| \leq M_1|u|^{-\nu}$ for $|u| > M_2$.
 194 A7 $h \rightarrow 0$, $nh \rightarrow \infty$ as $n \rightarrow \infty$.

195 Assumption A1 specializes the situation to the context of OU processes and has some
 196 relevant consequences for our results. First of all we note that any non-degenerate self-
 197 decomposable distribution is absolutely continuous (Sato 1999, Thm. 27.13). This, in turn,
 198 together with the postulated conditions on boundedness of the density and its derivatives (A2
 199 and A10 below) implies that the density $f = f_\theta$ belongs to the class of densities that satisfies

$$200 \sup_x f(x) + \sup_{x, x'} \frac{|f(x) - f(x')|}{|x - x'|} \leq M, \quad 0 < M < \infty \quad (6)$$

201 Second, from Masuda (2004, Theorem 4.3) it follows that $\{X_i\}_{0 \leq i \leq n}$ is ergodic and β -mixing
 202 with coefficients, for some $a > 0$, $\beta_X(t) = O(e^{-at})$. Recall that if X is a strictly stationary
 203 Markov process with initial distribution π and t^{th} step transition probability $P^t(x, \cdot)$, then the
 204 β -mixing coefficients are defined as

$$205 \beta_X(t) = \int \|P^t(x, \cdot) - \pi(\cdot)\| \pi(dx)$$

206 where $\|\mu\|$ denotes the total variation norm of a signed measure μ . The fact that $\{X_i\}_{0 \leq i \leq n}$
 207 is α -mixing follows from the inequality $2\alpha(t) \leq \beta(t)$. Conditions A2 and A3 require that
 208 all densities involved are bounded and A4 introduces a very mild tail restriction. Conditions
 209 A5 and A7 are quite standard in kernel density estimation while A6, introduced in Hansen
 210 (2008), is satisfied by most kernels including the normal one.

211 Our proof of consistency has a very mild restriction on existence of moments (A1 and A4)
 212 and uses boundedness and continuity (but not differentiability) conditions on the densities
 213 involved (A2, A3). On the other hand, it will require that the density estimates be restricted
 214 on a compact interval $\{x : |x| \leq c_n\}$ with $c_n \rightarrow \infty$ as $n \rightarrow \infty$ so that, ultimately, consistency
 215 will hold on a set of probability 1. The truncating device is defined in the [Appendix](#).

216 **Theorem 1.** Assume conditions A1–A7; then $|\hat{\theta}_1 - \theta_0| = o_p(1)$.

217 Asymptotic normality of $\hat{\theta}_1$ appears to need additional regularity assumptions, as well as
 218 existence of third order moments of X . This issue is investigate further in the simulations
 219 section.

220 **2.2. A minimum distance to independence estimator**

221 As we will see, a kernel based estimator suffers some problems when distributions with very
 222 heavy tails are involved. In such cases it may be sensible to resort to an alternative; here, in

231 order to provide an estimator which could be used under a minimal set of conditions and
 232 which could be a computationally attractive competitor, we introduce an estimator based on
 233 a minimum distance from independence. Define, with I_A being the indicator function of A ,

$$234 \hat{F}_\theta(t) = \frac{1}{n} \sum_{j=1}^n I_{(e_j^\theta \leq t)} \quad (7)$$

235 and

$$236 \hat{F}_\theta(t_1, t_2) = \frac{1}{n(n-1)} \sum_{i \neq j}^n I_{(e_i^\theta \leq t_1)} I_{(e_j^\theta \leq t_2)} \quad (8)$$

237 An estimator of θ can be obtained as

$$238 \hat{\theta}_2 = \arg \min_{\theta \in \Theta} \sup_{t_1, t_2 \in \mathbb{R}} \left| \hat{F}_\theta(t_1, t_2) - \hat{F}_\theta(t_1) \hat{F}_\theta(t_2) \right| \quad (9)$$

239 The use of the *sup* norm rather than other measures of distance is dictated by the desire to
 240 construct an estimator based on a minimal set of conditions on F . We then have (see Appendix
 241 for the proof).

242 **Theorem 2.** *Assume A1 then $|\hat{\theta}_2 - \theta_0| = o_p(1)$.*

243 In terms of computing $\hat{\theta}_2$, one may note that

$$244 \begin{aligned} 245 \hat{F}_\theta(t_1, t_2) - \hat{F}_\theta(t_1) \hat{F}_\theta(t_2) &= \frac{1}{n(n-1)} \sum_{i \neq j}^n I_{(e_j \leq t_1)} I_{(e_i \leq t_2)} - \frac{1}{n^2} \sum_{i,j=1}^n I_{(e_j \leq t_1)} I_{(e_i \leq t_2)} \\ 246 &= \frac{1}{n^2(n-1)} \sum_{i \neq j}^n I_{(e_j \leq t_1)} I_{(e_i \leq t_2)} - \frac{1}{n^2} \sum_{i=1}^n I_{(e_i \leq t_1)} I_{(e_i \leq t_2)} \end{aligned}$$

247 The actual computation of the estimators $\hat{\theta}_1$ and $\hat{\theta}_2$ can be done by a simple grid search.

248 3. Performance in finite samples

249 In this section the finite-sample performance of the proposed estimators is analyzed by
 250 simulations. The base-line to which we will compare the performance of our estimators will be
 251 the AC based estimator which is equivalent to several approaches proposed in the literature
 252 (see the introductory section for discussion about this) and, for the case of processes with
 253 positive increments, with the highly efficient estimator $\hat{\theta} = \min_{1 \leq j \leq n} X_j / X_{j-1}$ proposed by
 254 Jongbloed, Van der Meulen, and Van der Waart (2005); it is expected that $\hat{\theta}_1$ and $\hat{\theta}_2$ will
 255 not perform better than $\hat{\theta}$ however it is of interest here to give an overall evaluation of their
 256 performance.

257 Distributions over the real line such as the normal, the normal inverse Gaussian, the
 258 t -Student and the stable are considered; inverse Gaussian and stable OU processes with
 259 positive increments will also be used. The notation used will be a standard one, i.e., a normal
 260 distribution with mean μ and variance σ^2 will be denoted as $N(\mu, \sigma^2)$; the normal inverse

277 Gaussian distributions is indicated with $NIG(\alpha, \beta, \mu, \sigma)$ where α, β, μ and σ are related,
 278 respectively, to the tail, asymmetry, location and scale, $0 \leq \beta \leq \alpha, \mu \in \mathbf{R}, \sigma > 0$; t_ν stands
 279 for a t -Student distribution with ν degrees of freedom while $S(\alpha, \beta, \mu, \sigma)$ denotes a stable
 280 distribution with index of stability α , and where β, μ and σ indicate, respectively, asymmetry,
 281 location and scale; here we have $0 < \alpha \leq 2, 0 \leq \beta \leq 1, \mu \in \mathbf{R}, \sigma > 0$; the inverse Gaussian
 282 distribution with mean μ and shape σ will be indicated by $IG(\mu, \sigma)$.

283 As to the choice of kernel, we will compare two possibilities: a normal kernel, which is a
 284 standard choice in many computer packages, as well as a heavy tail kernel which should work
 285 better for heavy tailed distributions, namely

$$286 \quad 287 \quad 288 \quad 289 \quad K(u) = \frac{1}{2} e^{-|u|} \quad (10)$$

290 The choice of the smoothing bandwidth h exhibits a strong influence on the resulting estimate
 291 and it may not be optimal to consider automatic choices in running extensive simulations.
 292 Several alternatives have been compared: Silverman's rule of thumb, least squares cross
 293 validation, Sheather-Jones, over-smooth rule, standard deviation; we found that the choice of
 294 simply using the standard deviation as bandwidth works generally quite well for our problem
 295 and here we report estimation results based on that choice without any changes on single
 296 cases; this will allow a fair comparison on the estimators.

297 We found that the kernel-based estimators suffer some problems when facing distributions
 298 with heavy tails, where it is clear that in some cases the estimation procedure is failing
 299 completely, e.g. illogical results or improper kernel estimates. Hence implementation of
 300 formula Eq. (5) was carried out by eliminating those data for which $e > M$ for a given M .
 301 In the tables, simulated results with trimming and without trimming are reported; the value
 302 of M is indicated in the tables by writing $e \leq M$, i.e. all values $e > M$ have been eliminated.
 303 The choice of M is the result of a trial and error procedure by which the problems noted
 304 above are eliminated. For non-stable distributions no trimming was used. In the simulations,
 305 to prevent any bias in the comparison of the estimators we have chosen a general rule for
 306 trimming outliers and report the results as they are; we suspect that considering data-driven
 307 techniques would improve substantially the performance of the kernel-based method in the
 308 case of heavy-tailed distributions.

309 The OU processes with given marginal distribution have been generated according to the
 310 technique suggested in Taufer and Leonenko (2009b). All simulations have been run using
 311 the Mathematica® 8 software and the commands there automatically defined for kernel
 312 density estimation ("Smooth Kernel Distribution" with the "Standard deviation" bandwidth
 313 selection method) as well as for random number generation. The grid search for the value
 314 of θ maximizing the estimated likelihood or minimizing the Kolmogorov distance from
 315 independence has been set from 0.01 to 0.99 with 0.01 increments.

316 Tables 1–7 respectively report the estimation results for OU processes with marginal
 317 distributions: 1) $N(0, 3)$; 2) $NIG(2, 1.7, -1, 1)$; 3) t_4 ; 4) $S(1.5, -0.8, 0, 1)$; 5) $S(0.5, 0.5, 0, 1)$;
 318 6) $S(0.8, 1, 0, 1)$; 7) $IG(2, 2)$, the last two cases being positive distributions. For all cases but
 319 the IG one, where $\lambda = 0.5$, λ has been set to one. The examples proposed cover a variety of
 320 cases with symmetric and asymmetric, heavy and semi-heavy tailed marginal distributions.
 321 The Monte Carlo estimates of the mean and mean squared error (\widehat{MSE}) of the estimators of
 322 $\theta = e^{-\lambda}$ are based on 1000 simulations of samples with sizes $n = 50, 100, 200, 300$, where,

323 **Table 1.** Monte Carlo simulation results: $N(0,3); \theta = 0.3679$ ($\lambda = 1$). Mean, MSE and Relative efficiency (RE)
 324 of the estimators with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
326 AC	Mean	0.3236	0.3445	0.3560	0.3591
	MSE	0.0191	0.0953	0.0047	0.0027
328 SE	Mean	0.4476	0.3969	0.3784	0.3732
	RE	0.3086	0.6677	0.5550	0.5746
330 NO	Mean	0.3312	0.3476	0.3577	0.3603
	RE	1.0240	0.9777	0.9857	0.9958
332 HT	Mean	0.3297	0.3473	0.3576	0.3602
	RE	0.9822	0.9295	0.9487	0.9790

334 **Table 2.** Monte Carlo simulation results: $NIG(2,1.7,-1,1); \theta = 0.3679$ ($\lambda = 1$). Mean, MSE and Relative
 335 efficiency (RE) of the estimators with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
337 AC	Mean	0.3188	0.3403	0.3544	0.3588
	MSE	0.0170	0.0087	0.0041	0.0028
339 SE	Mean	0.4120	0.3759	0.3721	0.3694
	RE	0.5532	1.0148	1.1043	1.2807
341 NO	Mean	0.3319	0.3511	0.3619	0.3642
	RE	1.1018	1.3368	1.9082	2.0078
343 HT	Mean	0.3194	0.3590	0.3638	0.3632
	RE	1.1059	2.6348	2.7717	2.5150

345 **Table 3.** Monte Carlo simulation results: $t_4; \theta = 0.3679$ ($\lambda = 1$). Mean, MSE and Relative efficiency (RE) of
 346 the estimators with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
349 AC	Mean	0.3171	0.3375	0.3531	0.3590
	MSE	0.0180	0.0097	0.0045	0.0031
351 SE	Mean	0.4235	0.3832	0.3717	0.3670
	RE	0.3757	0.5184	0.6946	0.7699
353 NO	Mean	0.3302	0.3451	0.3568	0.3611
	RE	1.2514	1.2296	1.2833	1.2597
355 HT	Mean	0.3301	0.3462	0.3576	0.3617
	RE	1.2134	1.2472	1.3492	1.2905

357
 358 for an estimator $\hat{\theta}$:

$$359 \widehat{MSE}(\hat{\theta}) = \frac{1}{M} \sum_{i=1}^M (\hat{\theta}_i - \theta_0)^2 \quad (11)$$

360 with $\hat{\theta}_i$ the estimator obtained at the i -th Monte Carlo replicate, $i = 1, 2, \dots, M$.

361 Each table reports: mean and \widehat{MSE} for the auto-correlation estimator (AC); mean and
 362 relative efficiency (RE) with respect to the AC estimator for:

- 363 1. the minimum distance to independence estimator $\hat{\theta}_2$, indicated with SE;
- 364 2. the normal kernel-based $\hat{\theta}_1$ estimator, indicated with NO;
- 365 3. the Eq. (10) kernel-based $\hat{\theta}_1$ estimator, indicated with HT;

369 **Table 4.** Monte Carlo simulation results: Stable(1.5,−0.8,0,1); $\theta = 0.3679$ ($\lambda = 1$). Mean, MSE and Relative
 370 efficiency (RE) of the estimators with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
AC	Mean	0.3156	0.3427	0.3565	0.3594
	MSE	0.0152	0.0071	0.0032	0.0023
SE	Mean	0.4341	0.3941	0.3783	0.3735
	RE	0.3864	0.5180	0.6030	0.7908
NO	Mean	0.3398	0.3561	0.3646	0.4271
	RE	1.5493	1.7451	1.3533	0.0660
HT	Mean	0.3436	0.3592	0.3660	0.4279
	RE	1.6244	1.8969	1.4965	0.0670
HTC $e \leq 50$	Mean	0.3435	0.3592	0.3642	0.3670
	RE	1.6141	1.8171	1.8012	1.4878

382 **Table 5.** Monte Carlo simulation results: Stable(0.5, 0.5, 0,1); $\theta = 0.3679$ ($\lambda = 1$). Mean, MSE and Relative
 383 efficiency (RE) of the estimator with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
AC	Mean	0.3276	0.3462	0.3569	0.3622
	MSE	0.0085	0.0042	0.0024	0.0004
SE	Mean	0.3920	0.3763	0.3720	0.3706
	RE	1.4754	4.6537	14.960	16.7523
NO	Mean	0.3568	0.3598	0.4931	–
	RE	2.9792	2.9395	0.0383	–
HT	Mean	0.3704	0.3698	0.4977	–
	RE	4.9750	3.2251	0.0384	–
HTC $e \leq 500$	Mean	0.3788	0.3701	0.3652	0.3687
	RE	0.6538	0.7316	0.7937	0.5607

395 **Table 6.** Monte Carlo simulation results: S(0.8,1,0,1); $\theta = 0.3679$ ($\lambda = 1$). Mean, MSE and Relative efficiency
 396 (RE) of the estimators with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
AC	Mean	0.3238	0.3476	0.3582	0.3611
	MSE	0.0096	0.0038	0.0017	0.0009
SE	Mean	0.3855	0.3759	0.3725	0.3704
	RE	1.5854	3.3078	4.6422	4.7578
NO	Mean	0.3545	0.3625	0.4432	–
	RE	2.8035	3.6082	0.0480	–
HT	Mean	0.3683	0.3707	0.4459	–
	RE	6.3683	15.2199	0.0485	–
HTC $e \leq 100$	Mean	0.3654	0.3636	0.3640	0.3684
	RE	3.7733	2.8597	2.0772	2.1075
RA	Mean	0.4025	0.3851	0.3761	0.3728
	RE	4.0470	5.9336	10.3867	15.7843

411 4. the ratio-estimator of Jongbloed, Van der Meulen, and Van der Waart (2005) is indicated
 412 with RA.

413 In the case of stable distributions for which, as mentioned, automatic simulations with stan-
 414 dard settings suffered some problems, results for the kernel-based estimator HT computed

415 **Table 7.** Monte Carlo simulation results: $IG(2,2)$; $\theta = 0.6065$ ($\lambda = 0.5$). Mean, MSE and relative efficiency
 416 (RE) of the estimators with respect to AC. Estimates based on 1000 replications.

		$n = 50$	$n = 100$	$n = 200$	$n = 300$
418 AC	Mean	0.5491	0.5763	0.5927	0.5953
	MSE	0.0129	0.0058	0.0028	0.0020
420 SE	Mean	0.6493	0.6298	0.6151	0.6077
	RE	0.6597	0.6625	0.9352	0.9704
422 NO	Mean	0.5759	0.5904	0.6001	0.6008
	RE	1.9786	1.9568	1.8757	1.8876
424 HT	Mean	0.5881	0.5968	0.6034	0.6036
	RE	3.5109	3.8107	3.8725	4.0709
426 RA	Mean	0.6327	0.6282	0.6238	0.6222
	RE	15.9005	10.9213	8.6006	7.3762

427
428
429 with extremes outliers censored out are reported; this is indicated as *HTC* and the level M
 430 above which residuals have been eliminated is indicated as $e \leq M$.

431 The choice of reporting the RE with respect to the AC estimator is in order to emphasize the
 432 comparisons with respect to a cornerstone for all estimators. If $\hat{\theta}_{AC}$ denotes the AC estimator
 433 and $\hat{\theta}_O$ denotes any other estimator used in the simulations, then
 434

$$435 \quad RE(\hat{\theta}_O) = \frac{\widehat{MSE}(\hat{\theta}_{AC})}{\widehat{MSE}(\hat{\theta}_O)} \quad (12)$$

436
437
438 An RE higher than one results in a better performance of the estimator under analysis with
 439 respect to the AC estimator.
 440

441 In terms of investigating whether these estimators are asymptotically normal, Figures 1–4
 442 show the distribution of the estimators for some of the cases discussed in the tables, namely
 443 we consider the OU processes with $N(0, 3)$ and $IG(2, 2)$ marginal distribution either where λ
 444 is estimated using the normal kernel or the heavy kernel Eq. (10). In each figure the histogram
 445 of the standardized data is super-imposed with the standard normal density and PP and QQ
 446 plots for normality are reported. As we note from the figures, a normal approximation works
 447 quite well in all cases for sample sizes of around $n = 100$.

448 To summarize, the results in the tables and the figures are quite clear and indicate that
 449 generally the *NO* and *HT* estimators perform better with respect to the *AC* estimator having
 450 some problems only in the case of extremely heavy tails where in this case the *SE* estimator
 451 performs very well. Specifically we can summarize the results as follows:

- 452 a) in the Normal case the relative efficiency of *NO* and *HT* is always quite close to unity
 453 essentially indicating (as suggested by the theoretical results) no loss in efficiency with
 454 respect to the maximum likelihood estimator AC, even for small sample sizes.
 455 b) In all other cases the performance of *NO* and *HT* is generally better with respect to AC
 456 and relative efficiency can be quite high. In the stable case, some distinction needs to be
 457 made: it appears that, as sample size increases, the large number of extreme observations
 458 has a serious effect on the efficiency of the estimators, trimming can improve the
 459 situation. The tables, reporting the results of standardized simulations, may not show
 460 the effective performance of *NO* and *HT* in these cases.

461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506

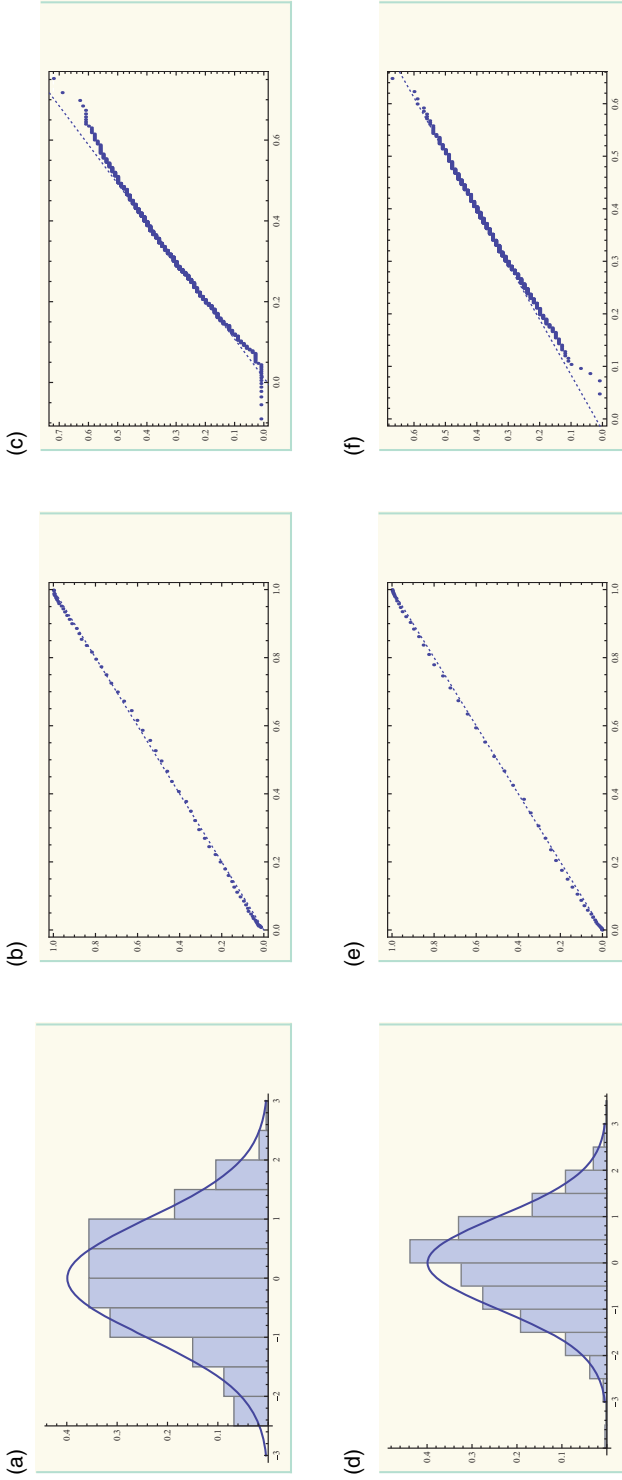


Figure 1. OU process with $\mathcal{N}(0, 3)$ marginal distribution. Empirical summaries of 1000 estimates of $\hat{\theta}_1$ based on the normal kernel; $n = 50$ (top) and $n = 100$ (bottom) sample size.

507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552

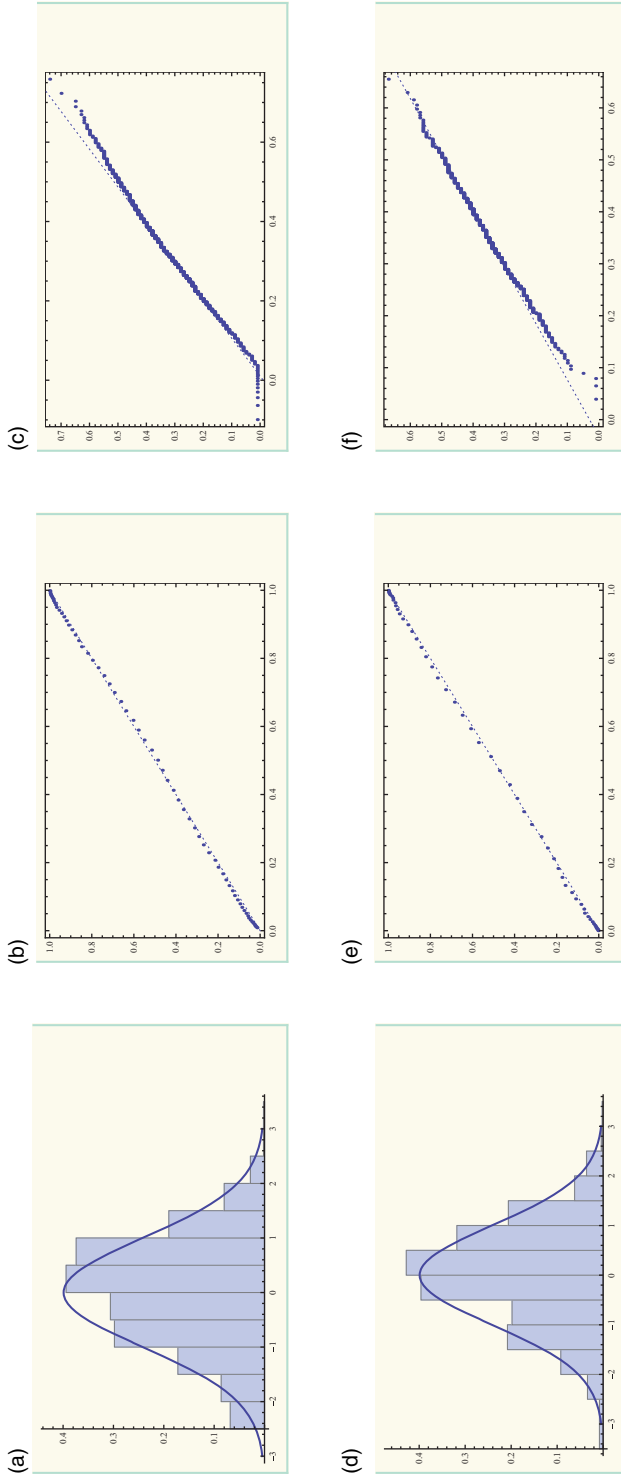


Figure 2. OU process with $N(0, 3)$ marginal distribution. Empirical summaries of 1000 estimates of $\hat{\theta}_1$ based on the heavy kernel Eq. (10); $n = 50$ (top) and $n = 100$ (bottom) sample size.

553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598

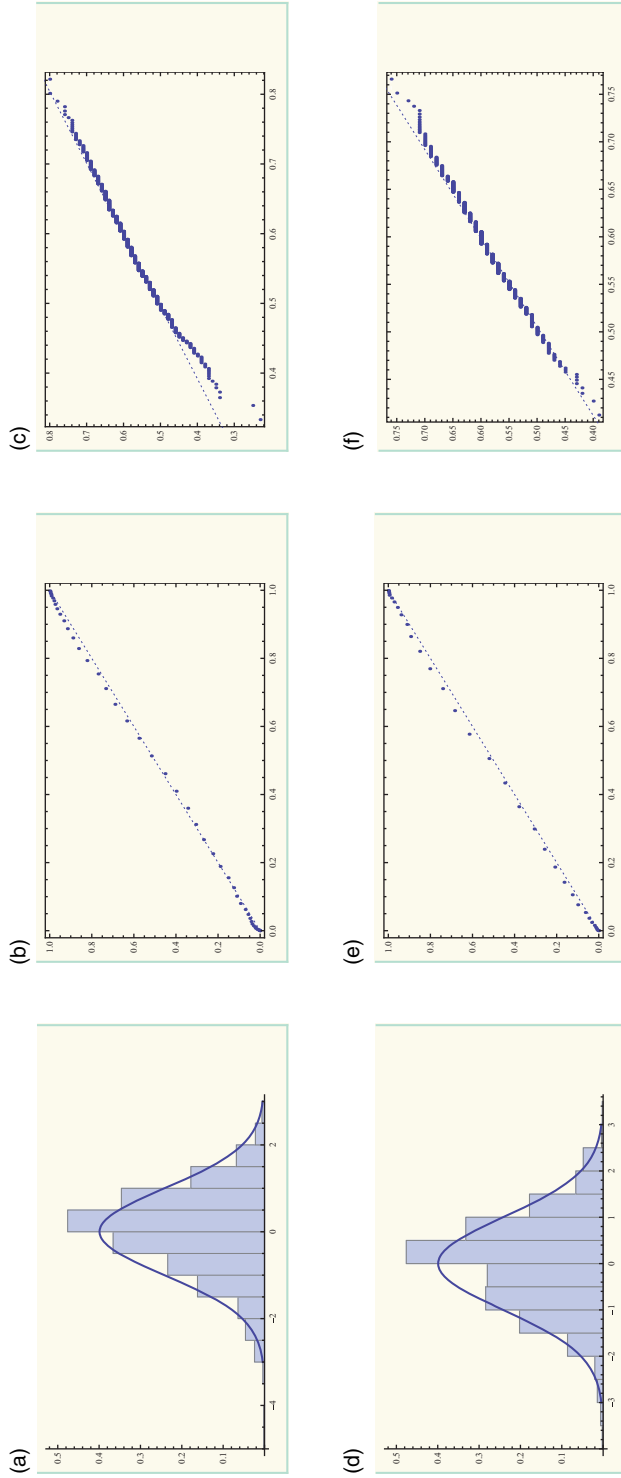


Figure 3. OU process with $/G(2, 2)$ marginal distribution. Empirical summaries of 1000 estimates of θ_1 based on the normal kernel; $n = 50$ (top) and $n = 100$ (bottom) sample size.

599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644

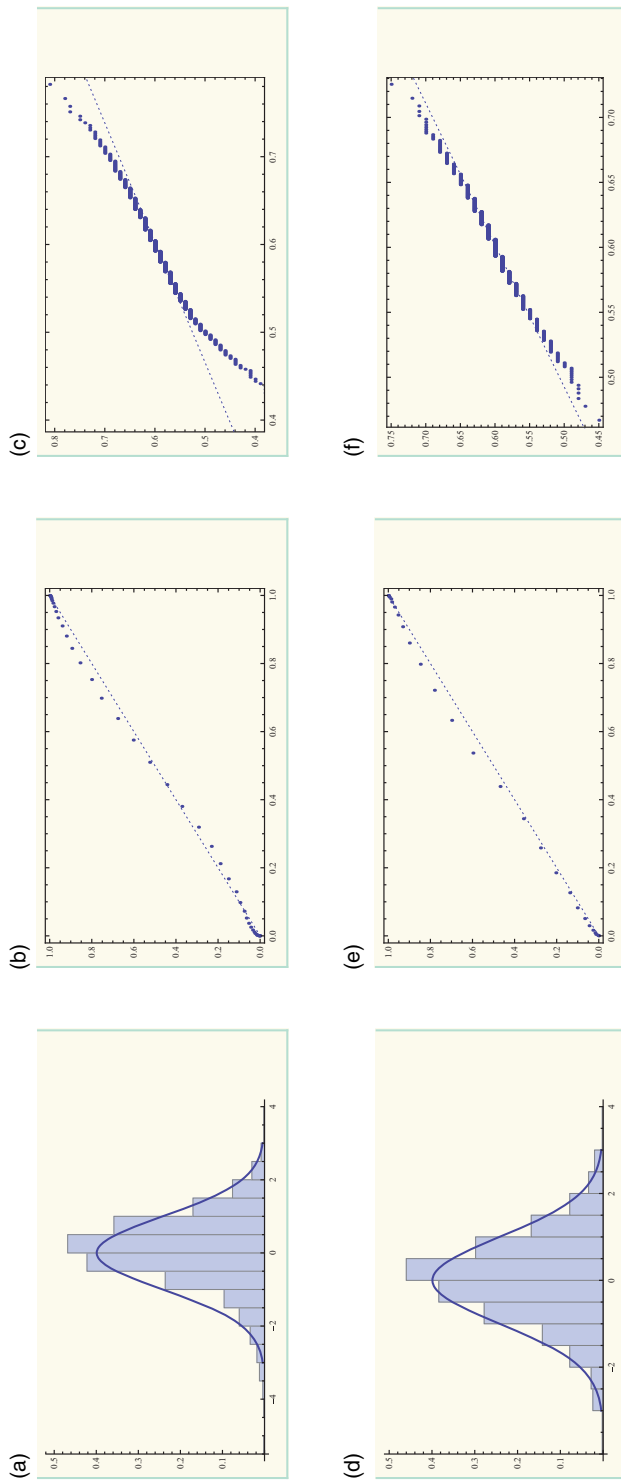


Figure 4. OU process with $G(2, 2)$ marginal distribution. Empirical summaries of 1000 estimates of $\hat{\theta}_1$ based on the heavy kernel Eq. (10); $n = 50$ (top) and $n = 100$ (bottom) sample size.

- 645 c) The performance of *HT* is generally better than *NO* and its relative efficiency can be much
 646 higher in semi-heavy or heavy tail cases.
- 647 d) In the case of OU processes with positive increments, the performance of *RA* is generally
 648 better than all the other estimators and its efficiency can be substantially larger than one.
 649 Note however that the *HT* estimator can perform extremely well for small sample sizes
 650 in the stable case and overcome the performance of *RA*.
- 651 e) The performance of *SE* is generally poorer with respect to the other estimators but in the
 652 case of distributions with very heavy tails, e.g. stable with $\alpha < 1$, for which *SE* does not
 653 suffer from the presence of extremely large observations.
- 654 f) A general rule, which seems to be efficient in a large variety of cases is the following: use
 655 *SE* if very heavy tails are present otherwise use *HT*. The *RA* estimator should be used for
 656 OU processes with positive increments.
- 657 g) Asymptotic normality seems to hold very well in all cases discussed.

660 4. An example

661
 662 Moody's trailing 12-month default rates are widely monitored indicators of corporate credit
 663 quality and are a good source either for theoretical and empirical studies. For example,
 664 Amerio, Muliere, and Secchi (2004) have studied the historical distributions of one-year
 665 default rates for Ba-rated, B-rated and Caa-rated defaulters during the period 1970–1999;
 666 Keenan, Sobehart, and Hamilton (1999) and Taufer (2007) have used either the entire
 667 Moody's rated universe (all-corporate, AC) and a sub-grouping, i.e., the speculative-grade
 668 (SG) monthly data respectively from 1970 to 1999 and from 1920 to 2004 in order to provide
 669 forecasting models.

670 In this example we are going to consider the SG yearly data for the period 1920–2011 for
 671 a total of 92 observations ranging from a minimum value of 0 to a maximum of 15.641. The
 672 data are taken from Moody's website and are freely available.

673 To begin with, we have a look at the linear plots of the series in Figure 5(a). The path
 674 does not appear to be non-stationary, however the high spikes suggests non normality of
 675 the data, which is confirmed by analytical tests and normality plot (not shown here). The
 676 auto-correlation and partial auto-correlation function in Figures 1(b) and 1(c) suggests that
 677 a (discretely observed) OU model could be appropriate for this data.

678 If normality is excluded, using the AC estimator maybe inappropriate; instead, one could
 679 consider some alternative approaches. Following the results of the simulations, for positive
 680 distributions, the highly efficient ratio estimator (*RA*) of Jongbloed, Van der Meulen, and
 681 Van der Waart (2005) should be used. Note however that the presence of several null values
 682 in the data makes it impossible its calculation. Also, the minimum distance estimator (*SE*)
 683 seems inappropriate here as there is no evidence of heavy tails and it is generally less efficient
 684 with respect to the kernel ones. Then, following the recommendations given in Section 3 we
 685 proceed to compute the *HT* estimator with no trimming. In this case the AC and the *HT*
 686 estimator are in good agreement, giving an estimate of 0.64 and 0.63 respectively. A further
 687 estimation on subsets of the data with 12 series of length 80 give values of the AC and *HT*
 688 estimators within a range of 0.63 and 0.71. Even though in this example the two estimators
 689 are very close, using alternative approaches is important in order to substantiate our empirical
 690 analysis.

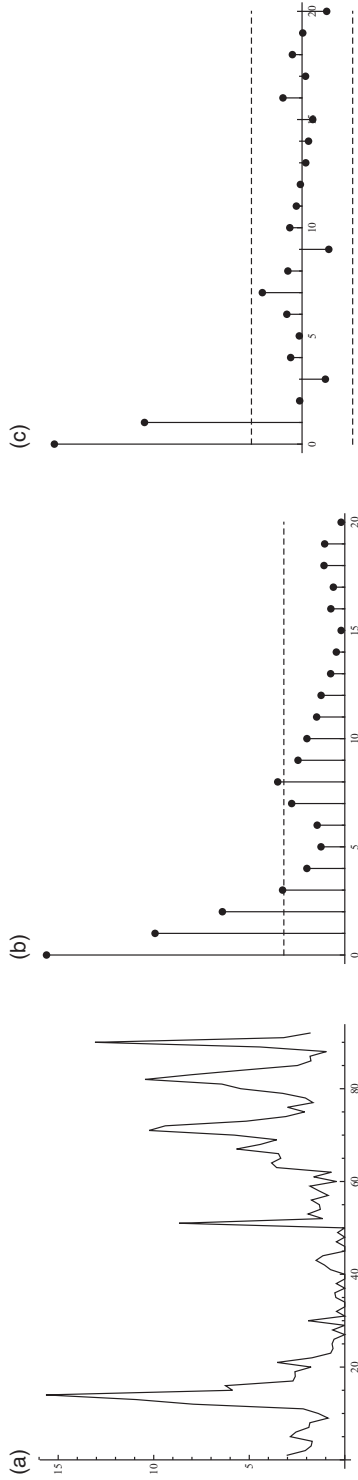


Figure 5. Autocorrelation and partial auto-correlation functions of the data.

691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736

References

- 737 Amerio, E., P. Muliere, and P. Secchi. 2004. Reinforced urn processes for modeling credit default
738 distributions. *International Journal of Theoretical and Applied Finance* 7 (4):407–23.
- 739 Andrews, B., M. Calder, and R. A. Davis. 2009. Maximum likelihood estimation for α -stable autore-
740 gressive processes. *The Annals of Statistics* 37:1946–82.
- 741 Barndorff-Nielsen, O. E. 1998. Processes of normal inverse Gaussian type. *Finance and Stochastics* 2:41–
742 68.
- 743 Barndorff-Nielsen, O. E., and N. Shephard. 2001. Non-Gaussian Ornstein–Uhlenbeck-based models
744 and some of their uses in financial economics. *Journal of the Royal Statistical Society, Series B.*
745 *Statistical Methodology* 63:167–241.
- 746 Brockwell, P. J., R. A. Davis, and Y. Yang. 2007. Estimation for nonnegative Lévy-driven Ornstein-
747 Uhlenbeck processes. *Journal of Applied Probability* 44:977–89.
- 748 Brown, D. J., and M. H. Wegkamp. 2002. Weighted minimum mean \hat{A} -square distance from indepen-
749 dence estimation. *Econometrica* 70:2035–51.
- 750 Chen, X., O. Linton, and I. Van Keilegom. 2003. Estimation of semiparametric models when the
751 criterion function is not smooth. *Econometrica* 71:1591–1608.
- 752 Diop, A., and A. F. Yode. 2010. Minimum distance parameter estimation for Ornstein–Uhlenbeck
753 processes driven by Lévy processes. *Statistics & Probability Letters* 80:122–7.
- 754 Drost, F. C., C. A. J. Klaassen, and B. J. M. Werker. 1997. Adaptive estimation in time-series models.
755 *Annals of Statistics* 25:786–817.
- 756 Durrett, R. A. 1991. *Probability: theory and examples*. Belmont, CA: Wadsworth.
- 757 Hall, P. 1986. On powerful distributional tests based on sample spacings. *Journal of Multivariate Analysis*
758 19:201–24.
- 759 Hall, P., and S. C. Morton. 1993. On the estimation of entropy. *Annals of the Institute of Statistical*
760 *Mathematics* 45:69–88.
- 761 Hansen, B. E. 2008. Uniform convergence rates for kernel estimation with dependent data. *Econometric*
762 *Theory* 24:726–48.
- 763 Heyde, C. C. and N. N. Leonenko. 2005. Student processes. *Advances in Applied Probability* 37 (2):342–
764 65.
- 765 Hu, Y., and H. Long. 2009. Least squares estimator for Ornstein-Uhlenbeck processes driven by α -stable
766 motions. *Stochastic Processes and their Applications* 119:2465–480.
- 767 Jongbloed, G., F. H., Van der Meulen, and A. Van der Waart. 2005. Nonparametric inference for Lévy-
768 driven Ornstein Uhlenbeck processes. *Bernoulli* 11:759–91.
- 769 Kapur, J. N., and H. K. Kesavan. 1992. *Entropy optimization principles with applications*. Boston, MA:
770 Academic Press.
- 771 Keenan, S. C., J. Sobehart, and D. T. Hamilton. 1999. Predicting default rates: a forecasting model for
772 Moody’s issuer-based default rates. Moody’s Investor Service: Special Comment.
- 773 Koul, H. L. 1986. Minimum distance estimation and goodness-of-fit tests in first-order auto-regression.
774 *Annals of Statistics* 14:1194–213.
- 775 Koul, H. L., and A. Schick. 1997. Efficient estimation in nonlinear autoregressive time-series models.
776 *Bernoulli* 3:247–77.
- 777 Kreiss, J. P. 1987. On adaptive estimation in stationary ARMA processes. *Annals of Statistics* 15 (1):112–
778 33.
- 779 Linton, O., and Z. Xiao. 2007. A nonparametric regression estimator that adapts to error distribution
780 of unknown form. *Econometric Theory* 23 (3):371–413.
- 781 Linton, O., S. Sperlich, and I. Van Keilegom. 2008. Estimation of a semiparametric transformation
782 model. *Annals of Statistics* 36 (2):686–718.
- 783 Long, H. 2009. Least squares estimator for discretely observed Ornstein–Uhlenbeck processes with
784 small Lévy noises. *Statistics & Probability Letters* 79:2076–85.
- 785 Ma, C. 2010. A note on “Least squares estimator for discretely observed Ornstein-Uhlenbeck processes
786 with small Lévy noises”. *Statistics & Probability Letters* 80:1528–31.
- 787 Manski, C. F. 1983. Closest empirical distribution estimation. *Econometrica* 51:305–19.

- 783 Masuda, H. 2004. On multidimensional Ornstein–Uhlenbeck processes driven by a general Lévy
784 process. *Bernoulli* 10 (1):97–120.
- 785 Meintanis, S. G., and E. Taufer. 2012. Inference procedures for stable-Paretian stochastic volatility
786 models. *Mathematical and Computer Modelling* 55 (3–4):1199–212.
- 787 Newey, W. K. 1991. Uniform convergence in probability and stochastic equicontinuity. *Econometrica*
788 59:1161–67.
- 789 Nielsen, B., and N. Shephard. 2003. Likelihood analysis of a first-order autoregressive model with
790 exponential innovations. *Journal of Time Series Analysis* 24 (3):337–44.
- 791 Sato, K. 1999. *Lévy processes and infinitely divisible distributions*. Cambridge: Cambridge University
792 Press.
- 793 Schoutens, W., and J. L. Teugels. 1998. Lévy processes, polynomials and martingales. Special issue in
794 honor of Marcel F. Neuts. *Communications in Statistics: Stochastic Models* 14:335–49.
- 795 Seneta, E. 2004. Fitting of variance-Gamma model to financial data. *Journal of Applied Probability*
796 41:177–87.
- 797 Taufer, E. 2007. Modeling stylized features in default rates. *Applied Stochastic Models in Business and*
798 *Industry* 23:73–82.
- 799 Taufer, E., and N. Leonenko. 2009a. Characteristic function estimation of non-Gaussian Ornstein-
800 Uhlenbeck processes. *Journal of Statistical Planning and Inference* 139:3050–63.
- 801 Taufer, E., and N. Leonenko. 2009b. Simulation of Lévy-driven Ornstein-Uhlenbeck processes with
802 given marginal distribution. *Computational Statistics & Data Analysis* 53:2427–37.
- 803 Taufer, E., N. Leonenko, and M. Bee. 2011. Characteristic function estimation of Ornstein-Uhlenbeck-
804 based stochastic volatility models. *Computational Statistics & Data Analysis* 55:2525–39.
- 805 Van der Vaart, A. W. 1998. *Asymptotic statistics*. New York: Cambridge University Press.
- 806 van Es, B. 1992. Estimating functionals related to a density by a class of statistics based on spacings.
807 *Scandinavian Journal of Statistics* 19:61–72.
- 808 Yao, W., and Z. Zhao. 2013. Kernel density based linear regression estimate. *Communications in*
809 *Statistics—Theory and Methods* 42:4499–512.
- 810 Yuan, A., and J. De Gooijer. 2007. Semiparametric regression with kernel error model. *Scandinavian*
811 *Journal of Statistics* 34:841–69.
- 812 Zhang, S., and X. Zhang. 2013. A least squares estimator for discretely observed Ornstein–Uhlenbeck
813 processes driven by symmetric α -stable motions. *Annals of the Institute of Statistical Mathematics*
814 65:89–103.
- 815 Zhang, S., Z. Lin, and X. Zhang. 2015. A least squares estimator for Lévy-driven moving averages based
816 on discrete time observations. *Communications in Statistics—Theory and Methods* 44:1111–29.
- 817 Zhao, Z., and W. B. Wu. 2009. Nonparametric inference of discretely sampled stable Lévy processes.
818 *Journal of Econometrics* 153:83–92.

AQ5

814 Appendix

815 For the proof of the results, a preliminary lemma due to Hansen (2008), specialized to our set-up, is
816 needed.

817 **Lemma A1.** Assume conditions A1–A7; then, for $c_n = O((\ln n)n^{1/2})$,

$$818 \sup_{|x| \leq c_n} \left| \hat{f}_\theta(x) - f_\theta(x) \right| = O_p \left(\left(\frac{\ln n}{nh} \right)^{1/2} + h^q \right), \quad \forall \theta \in \Theta \quad (A1)$$

819 where q denotes the order of the kernel.

820 **Lemma A1** follows directly from Theorem 2 and Theorem 6 in Hansen (2008) by noting that, from
821 Masuda (2004), the sequence X_0, \dots, X_n following model Eq. (2) is β -mixing with geometric rate; we
822 therefore can take in the theorems of Hansen (2008), $\theta = 1$. The result of Lemma 1 can be strengthened
823 to almost sure convergence and convergence over the whole real line by strengthening the assumptions;
824 but the present version will suffice for our purposes.

AQ6

829 In the proof of the results, a smooth trimming function G_b will be used, where

830
831
$$G_b(x) = \begin{cases} 0, & x < b \\ \int_b^x g_b(z) dz, & b \leq x \leq 2b \\ 1, & x > 2b \end{cases} \tag{A2}$$

832
833 Here $g_b(x) = \frac{1}{b}g(x/b - 1)$ with $b > 0$ a trimming parameter and g any density function with support in $[0, 1]$, $g(0) = g(1) = 0$. This approach has been followed, for example, by Linton and Xiao (2007) and Yao and Zhao (2013) and a proper choice of g allows to use standard Taylor series arguments; for example, if $g(z) = cz^\alpha(1 - z)^\alpha$, $z \in [0, 1]$ $\alpha > 0$ and c an appropriate normalizing constant, then G_b is $\alpha + 1$ times continuously differentiable on $[0, 1]$. Note also that $\sup_x G_b(x)/x^k \leq 1/b^k$.

839 **Lemma A2.** Assume conditions A1–A7, then,
840 a)

841
$$\max_{1 \leq i \leq n} \left| \frac{\hat{f}_\theta(e_i) - f_\theta(e_i)}{f_\theta(e_i)} \right| = o_p(1), \quad \forall \theta \in \Theta \tag{A3}$$

843 b)

844
$$\sup_{|\theta_1 - \theta_2| \leq \varepsilon} \max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_1}(e_i) - \hat{f}_{\theta_2}(e_i)}{\hat{f}_{\theta_2}(e_i)} \right| = o_p(1) + O(\varepsilon) \tag{A4}$$

847 *Proof of Lemma A2.* For the proof of part a), consider first a trimmed version

848
849
$$\max_{1 \leq i \leq n} \left| \frac{\hat{f}_\theta(e_i) - f_\theta(e_i)}{f_\theta(e_i)} \right| G_b(f_\theta(e_i)) \leq \max_{1 \leq i \leq n} \frac{|\hat{f}_\theta(e_i) - f_\theta(e_i)|}{b} \tag{A5}$$

850 using the fact that $\sup_x G_b(x)/x \leq 1/b$. Next, for $\mathbf{I}_{(x)}$ the indicator function, note that

851
852
$$\max_{1 \leq i \leq n} |\hat{f}_\theta(e_i) - f_\theta(e_i)| \mathbf{I}_{(|e_i| \leq c_n)} \leq \sup_{|x| \leq c_n} |\hat{f}_\theta(x) - f_\theta(x)| \tag{A6}$$

853 hence Lemma A1, as $n \rightarrow \infty$, implies that Eq. (A5) is $o_p(1)O(1/b)$ and the result follows as the choice of b is arbitrary.

856 As far as part b) is concerned, using part a) we have,

857
858
$$\max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_2}(e_i)}{f_{\theta_2}(e_i)} - 1 \right| = o_p(1) \quad \text{and} \quad \max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_1}(e_i) - f_{\theta_1}(e_i)}{f_{\theta_2}(e_i)} \right| = o_p(1) \tag{A7}$$

859 Next note that

860
861
$$\frac{\hat{f}_{\theta_1}(e)}{\hat{f}_{\theta_2}(e)} = \frac{\hat{f}_{\theta_1}(e)/f_{\theta_2}(e)}{\hat{f}_{\theta_2}(e)/f_{\theta_2}(e)} = \frac{\frac{f_{\theta_1}(e)}{f_{\theta_2}(e)} + \frac{\hat{f}_{\theta_1}(e) - f_{\theta_1}(e)}{f_{\theta_2}(e)}}{1 + \frac{\hat{f}_{\theta_2}(e) - f_{\theta_2}(e)}{f_{\theta_2}(e)}} \tag{A8}$$

864 results in Eq. (A7) imply that

865
866
$$\max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_1}(e)}{\hat{f}_{\theta_2}(e)} - \frac{f_{\theta_1}(e)}{f_{\theta_2}(e)} \right| = o_p(1) \tag{A9}$$

868 Based on the above results we obtain

869
870
$$\begin{aligned} & \sup_{|\theta_1 - \theta_2| \leq \varepsilon} \max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_1}(e_i) - \hat{f}_{\theta_2}(e_i)}{\hat{f}_{\theta_2}(e_i)} \right| \\ & \leq \sup_{|\theta_1 - \theta_2| \leq \varepsilon} \max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_1}(e_i)}{\hat{f}_{\theta_2}(e_i)} - \frac{f_{\theta_1}(e_i)}{f_{\theta_2}(e_i)} \right| + \sup_{|\theta_1 - \theta_2| \leq \varepsilon} \max_{1 \leq i \leq n} \left| \frac{f_{\theta_1}(e_i)}{f_{\theta_2}(e_i)} - 1 \right| \\ & \leq o_p(1) + \frac{C\varepsilon}{b}; \quad \forall b \end{aligned} \tag{A10}$$

871
872
873
874

875 where the first term on the *r.h.s.* of the above expression is from Eq. (A8) whereas the second is again
 876 obtained by truncation and from condition A2. Again we can make the above term as small as desired
 877 as the choice on b and ε are arbitrary. \square

878 *Proof of Theorem 1.* In order to prove consistency of $\hat{\theta}_1$ we need to show that:

- 879 a) there is a function, say $L(\theta)$, such that $\sup_{\theta \in \Theta} |L_n(\theta) - L(\theta)| = o_p(1)$;
- 880 b) $L(\theta)$ is uniquely maximized by θ_0 .

881 In order to prove part *a*) we need to verify that: (i) the parameter space is compact; (ii) $L_n(\theta) \rightarrow_p$
 882 $L(\theta)$ point wise; (iii) equicontinuity in probability, i.e. there exists $\delta > 0$ such that $\sup_{|\theta_1 - \theta_2| \leq \delta} |L_n(\theta_1) -$
 883 $L_n(\theta_2)| = o_p(1)$.

884 As far as point (i) is concerned note that although $\Theta = (0, 1)$ is not compact we can consider a
 885 compact set K such that $\theta_0 \in K \subset (0, 1)$. In order to verify (ii), define $M_n = M_n(\theta) = \frac{1}{n} \sum_{j=1}^n \log f_{\theta}(e_j)$
 886 and $L(\theta) = E(\ln f_{\theta}(e))$. Then, since X_1, \dots, X_n is ergodic, under Assumption A4 it follows that $M_n \rightarrow_p$
 887 $L(\theta)$, $\forall \theta \in \Theta$. Since $|\ln(1+x)| \leq 2|x|$ in a neighborhood of $x = 0$, a sufficient condition for
 888 $L_n - M_n \rightarrow_p 0$, is

$$889 \max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta}(e_i)}{f_{\theta}(e_i)} - 1 \right| = o_p(1) \quad \forall \theta \in \Theta \tag{A11}$$

890 which follows from Lemma A2a. It follows that $L_n(\theta) \rightarrow_p L(\theta)$ point-wise. Similarly, to show *iii*) note
 891 that,

$$892 \sup_{|\theta_1 - \theta_2| \leq \varepsilon_n} \max_{1 \leq i \leq n} |L_n(\theta_1) - L_n(\theta_2)| = \sup_{|\theta_1 - \theta_2| \leq \varepsilon_n} \max_{1 \leq i \leq n} \frac{1}{n} \left| \sum_{i=1}^n \log \left(1 + \frac{\hat{f}_{\theta_1}(e_i^{\theta_1}) - \hat{f}_{\theta_2}(e_i^{\theta_2})}{\hat{f}_{\theta_2}(e_i^{\theta_2})} \right) \right|$$

$$893 \leq 2 \sup_{|\theta_1 - \theta_2| \leq \varepsilon_n} \max_{1 \leq i \leq n} \left| \frac{\hat{f}_{\theta_1}(e_i^{\theta_1}) - \hat{f}_{\theta_2}(e_i^{\theta_2})}{\hat{f}_{\theta_2}(e_i^{\theta_2})} \right|$$

894 which is $o_p(1)$ by Lemma A2b for suitably chosen ε_n . In order to prove part *b*), define $L(\theta) = -H(\varepsilon^{\theta})$
 895 where H is the Shannon's entropy (see, e.g., Kapur and Kesavan 1992). Then,

$$896 H(\varepsilon^{\theta}) = H(\varepsilon^{\theta_0} + (\theta_0 - \theta)X_0)$$

$$897 \geq H(\varepsilon^{\theta_0} + (\theta_0 - \theta)X_0|X_0)$$

$$898 = H(\varepsilon^{\theta_0}|X_0) \tag{A12}$$

$$899 = H(\varepsilon^{\theta_0})$$

900 where we have used, in order, the facts that; conditioning reduces entropy; a constant does not change
 901 entropy; ε^{θ_0} and X_0 are independent. It follows that $L(\theta)$ is uniquely maximized by $L(\theta_0)$. \square

902 *Proof of Theorem 2.* Denote for simplicity $\sup_{t_1, t_2 \in R} |\hat{F}_{\theta}(t_1, t_2) - \hat{F}_{\theta}(t_1)\hat{F}_{\theta}(t_2)| = \rho(\hat{F}, \theta)$. The proof of
 903 the theorem follows from Theorem 2 in Manski (1983) if we verify the following conditions:

- 904 B1 The parameter space Θ is compact.
- 905 B2 $\rho(F, \theta) = 0$ if and only if $\theta = \theta_0$.
- 906 B3 (Assumption 4 in Manski (1983) - continuity and uniform convergence). $\rho(F, \theta)$ is continuous as
 907 a function on Θ . Also, $\rho(\hat{F}, \theta)$ converges in probability to $\rho(F, \theta)$ uniformly over Θ .

908 As far as B1 is concerned, as already discussed, one can consider a compact set K such that $\theta_0 \in K \subset$
 909 $(0, 1)$. B2 follows from the discussion in Section 2, as the sequence $\{e_i^{\theta}\}_{1 \leq i \leq n}$ is i.i.d only if $\theta = \theta_0$.

910 The first part of B3 can be verified by first noting that F , being self-decomposable, is absolutely
 911 continuous (Sato 1999, Thm 27.13) and exploiting the first part of the corollary to Theorem 2 in Manski
 912 (1983) by noting that $g(X_1, X_0, \theta) = X_1 - \theta X_0$ is continuous on $S \times \Theta$ where $S \in R^2$ is some compact
 913 and convex set.

914 The second part follows if we prove that

$$915 \sup_{\theta \in \Theta} \sup_{t_1, t_2 \in R} \left| \hat{F}_{\theta}(t_1, t_2) - \hat{F}_{\theta}(t_1)\hat{F}_{\theta}(t_2) - F_{\theta}(t_1, t_2) + F_{\theta}(t_1)F_{\theta}(t_2) \right| = o_p(1) \tag{A13}$$

921 In order to do this, note that

$$\begin{aligned}
 & \left| \hat{F}_\theta(t_1, t_2) - \hat{F}_\theta(t_1)\hat{F}_\theta(t_2) - F_\theta(t_1, t_2) + F_\theta(t_1)F_\theta(t_2) \right| \\
 & \leq \left| \hat{F}_\theta(t_1, t_2) - F_\theta(t_1, t_2) \right| + \hat{F}_\theta(t_1) \left| \hat{F}_\theta(t_2) - F_\theta(t_2) \right| + F_\theta(t_2) \left| \hat{F}_\theta(t_1) - F_\theta(t_1) \right|
 \end{aligned} \tag{A14}$$

926 Note that since the sequence $\{X_i\}_{0 \leq i \leq n}$ is ergodic and the class of functions $\mathcal{F} = \{f_t = \mathbf{I}_{(-\infty, t]}, t \in \mathbb{R}^2\}$
 927 are Glivenko–Cantelli (see, e.g. Van der Vaart 1998, p. 270), we have that

$$\sup_{t_1, t_2 \in \mathbb{R}^2} \left| \hat{F}_\theta(t_1, t_2) - F_\theta(t_1, t_2) \right| = o_p(1), \quad \text{and} \quad \sup_{t \in \mathbb{R}} \left| \hat{F}_\theta(t) - F_\theta(t) \right| = o_p(1) \quad \forall \theta \in \Theta \tag{A15}$$

931 We claim that

$$\sup_{\theta \in \Theta} \sup_{t \in \mathbb{R}} \left| \hat{F}_\theta(t) - F_\theta(t) \right| = o_p(1) \tag{A16}$$

$$\sup_{\theta \in \Theta} \sup_{t_1, t_2 \in \mathbb{R}} \left| \hat{F}_\theta(t_1, t_2) - F_\theta(t_1, t_2) \right| = o_p(1) \tag{A17}$$

936 The proof of Eqs. (A16) and (A17) together with compactness of Θ and Eq. (A15) will prove Eq. (A13).

937 In order to prove Eqs. (A16) and (A17) we'll exploit Theorem 3 in Chen, Linton, and Van Keilegom
 938 (2003) which provides primitive conditions for equicontinuity: we'll have to show that their condition
 939 (3.2) is satisfied, which require in our case to show that

$$\left[\mathbb{E} \left(\sup_{|\theta_1 - \theta_2| \leq \delta} \left| \mathbf{I}_{\{X_1 - \theta_1 X_0 \leq t\}} - \mathbf{I}_{\{X_1 - \theta_2 X_0 \leq t\}} \right|^r \right) \right]^{1/r} \leq K\delta^s \tag{A18}$$

$$\left[\mathbb{E} \left(\sup_{|\theta_1 - \theta_2| \leq \delta} \left| \mathbf{I}_{\{X_1 - \theta_1 X_0 \leq t\}} \mathbf{I}_{\{X_2 - \theta_1 X_1 \leq t\}} - \mathbf{I}_{\{X_1 - \theta_2 X_0 \leq t\}} \mathbf{I}_{\{X_2 - \theta_2 X_1 \leq t\}} \right|^r \right) \right]^{1/r} \leq K\delta^s \tag{A19}$$

946 for all $\theta \in \Theta$, all small positive values $\delta = o(1)$, $r \geq 2$ and $s \in (0, 1]$ and that the bounds hold
 947 for μ -almost all (t_1, t_2) . Consider Eq. (A18) and note that the expectation of the absolute value in the
 948 expression is the probability of the union of the events $\{t + \theta_1 X_0 < X_1 < t + \theta_2 X_0\}$ and $\{t + \theta_2 X_0 <$
 949 $X_1 < t + \theta_1 X_0\}$ which consider all possibilities arising from the cases $\theta_1 \leq \theta_2$ or $\theta_1 > \theta_2$, $X_0 \leq 0$ or
 950 $X_0 > 0$. Since X_0 is bounded in probability there is a compact set with probability greater than $1 - \varepsilon$,
 951 $\varepsilon > 0$, for which there is some upper bound c such that $\sup_{|\theta_1 - \theta_2| \leq \delta} |\theta_1 X_0 - \theta_2 X_0| \leq \delta c$. For some
 952 $\delta > 0$ we have then

$$\begin{aligned}
 \mathbb{E} \left(\sup_{|\theta_1 - \theta_2| \leq \delta} \left| \mathbf{I}_{\{X_1 - \theta_1 X_0 \leq t\}} - \mathbf{I}_{\{X_1 - \theta_2 X_0 \leq t\}} \right| \right) & \leq 2P(t - \delta c + \theta X_0 < X_1 < t + \delta c + \theta X_0) \\
 & = 2F_\theta(t - \delta c) - F_\theta(t + \delta c) \\
 & \leq K\delta
 \end{aligned} \tag{A20}$$

956 for some constant $K < \infty$, from continuity of F . Therefore condition (3.2) of Theorem 3 in Chen,
 957 Linton, and Van Keilegom (2003) is satisfied with $r = 2$ and $s = 1/2$. The proof of Eq. (A19) resorts to
 958 an analogous device. \square

959
960
961
962
963
964
965
966