

Online Supplement for “Bootstrap Prediction Inference of Non-linear Autoregressive Models”

Kejin Wu¹ and Dimitris N. Politis ^{*,2}

¹Department of Mathematics, University of California, San Diego

²Department of Mathematics and Halicioğlu Data Science Institute, University of California,
San Diego

The supplement contains three parts. One set of sufficient conditions for geometric ergodicity of

$$X_t = \phi(\mathbf{X}_{t-1}) + \sigma(\mathbf{X}_{t-1})\epsilon_t \quad (1)$$

are discussed in [Appendix A](#). The application of the forward bootstrap prediction algorithm on the general NLAR model

$$X_t = G(\mathbf{X}_{t-1}, \epsilon_t)$$

is presented in [Appendix B](#). The proofs for the lemmas and theorems in the main text are given in [Appendix C](#).

Appendix A SUFFICIENT CONDITIONS FOR GEOMETRIC ERGODICITY

Checking geometric ergodicity for a LAR model is simple; it is well known that the LAR model is stationary and geometrically ergodic as long as the corresponding characteristic polynomial does not have zero roots inside or on the unit circle. However, this check criterion depends on the linearity assumption, and can not be extended to serve for NLAR model directly An and Huang (1996). Thus, practitioners rely on Markov chain techniques to explore conditions under which the NLAR model is geometrically ergodic. The motivation is that the NLAR model can be described as a Markov chain in a general state space; the extensive discussions and literature related to Markov chains can guide the development of criteria to check the ergodicity of NLAR models.

One of the earliest criteria developed to guarantee the ergodicity of a Markov chain is Doeblin’s condition given by Doob (1953). Later, Tweedie (1975) proposed a more generalized condition, the so-called *Drift criterion*. This criterion gives a sufficient condition for an aperiodic and irreducible Markov chain to be geometrically ergodic. For completeness, we present this criterion using the version applied by An and Huang (1996):

*Correspondence to: Dimitris N. Politis. Email: dpolitis@ucsd.edu.

Lemma A.1 (Drift criterion). *Let $\{X_t\}$ be an aperiodic and irreducible Markov chain. Suppose that there exists a small set C , a non-negative measurable function k , positive constants c_1, c_2 and $\rho < 1$ such that:*

$$\begin{aligned} \mathbb{E}(k(\mathbf{X}_{t+1})|\mathbf{X}_t = \mathbf{x}) &\leq \rho k(\mathbf{x}) - c_1, \text{ for any } \mathbf{x} \notin C, \\ \mathbb{E}(k(\mathbf{X}_{t+1})|\mathbf{X}_t = \mathbf{x}) &\leq c_2, \text{ for any } \mathbf{x} \in C. \end{aligned} \quad (2)$$

Remark. *The function $k(\cdot)$ is called the test function in the literature. For the formal definition of a ‘small set’; see Tjøstheim (1990) or Tong (1990). We will soon see that the small set can be taken as a compact set in some situations.*

Thus, to ensure the ergodicity of an NLAR model, people can check if $\{X_t\}$ is aperiodic and irreducible and Lemma A.1 holds. Along with this idea, An and Huang (1996) give several kinds of sufficient conditions for Eq. (1) with $\sigma(\mathbf{X}_{t-1}) \equiv 1$ (homoscedastic errors case) to be geometrically ergodic. Note that the test function $k(\cdot)$ and the specific small set will change according to which condition is based. Then, Min and Hongzhi (1999) extend these results to the region of NLAR models with heteroscedastic errors. Based on this body of work, if we assume:

A4 The probability density function of innovation $f_\epsilon(\cdot)$ is continuous and everywhere positive.

A5 The conditional mean and volatility functions satisfy the inequalities:

$$\sup_{\|\mathbf{x}\|_2 \leq K} |\phi(\mathbf{x})| < \infty; \quad \sup_{\|\mathbf{x}\|_2 \leq K} |\sigma(\mathbf{x})| < \infty, \text{ for each } K > 0, \quad (3)$$

where $\mathbf{x} \in \mathbb{R}^p$, and $\|\cdot\|_2$ is the Euclidean norm.

we can obtain a useful Lemma:

Lemma A.2. *Let $\{X_t\}$ satisfy the model Eq. (1). Suppose A4 and A5 are fulfilled. Then $\{X_t\}$ is aperiodic and irreducible with respect to μ which is the Lebesgue measure. Moreover, the μ -non-null compact sets are small sets.*

Remark A.1. *The proof of Lemma A.2 can be found in Min and Hongzhi (1999). The original proof only requires the density function to be lower semi-continuous. In this paper, since the estimation of the innovation distribution and the L_1 optimal predictor will be discussed, we require a stronger condition for the density function. Besides, we should notice that we can check the classical properties defined by the Markov chain to verify the aperiodicity and irreducibility. Nevertheless, it is more natural to apply Lemma A.2 for analyzing NLAR models. For example, when we consider a homoscedastic NLAR model with $p = 1$, the drift criterion and the boundedness of the conditional mean function on a suitable compact set can be satisfied by a simple condition: $|\phi(x)| \leq C_2|x| + C_1$, for all x and some $C_1 < \infty, C_2 < 1$; this is exactly Assumption 3 (i) of Franke, Kreiss, et al. (2002). In their work, they mentioned that the everywhere-positive assumption (A4) on the density function is unnecessarily restrictive. Thus, they replace it with their Assumption 3 (iii). However, for simplifying the proofs of our paper, we still require the everywhere-positive property. Furthermore, we will apply a convolutional technique to acquire an everywhere-positive “innovation distribution” in the bootstrap world to satisfy this strong assumption.*

The condition which ensures the Drift criterion is satisfied for a homoscedastic NLAR models can be drawn from An and Huang (1996). If we assume:

A6 There exists a positive number $\lambda < 1$ and a constant C such that the conditional mean function satisfies:

$$|\phi(\mathbf{x})| \leq \lambda \max\{|x_1|, \dots, |x_p|\} + C, \quad (4)$$

then we can get the below lemma:

Lemma A.3 (Theorem 3.2 of An and Huang (1996)). *Let $\{X_t\}$ satisfy the homoscedastic version of the model Eq. (1). If A4 and A6 are fulfilled, then this NLAR model is geometrically ergodic.*

Similarly, we can obtain the condition for heteroscedastic NLAR models to be geometrically ergodic by imposing an additional assumption on the variance function:

A7 The conditional variance function satisfies:

$$\lim_{\|\mathbf{x}\|_2 \rightarrow \infty} \frac{\sigma(\mathbf{x})}{\|\mathbf{x}\|_2} = 0. \quad (5)$$

Then the Lemma holds:

Lemma A.4 (Theorem 3.5 of Min and Hongzhi (1999)). *Let $\{X_t\}$ satisfy the model Eq. (1). Suppose the conditional variance function satisfies A5 and A7. In addition, A6 holds true for the conditional mean function, and A4 holds true for the probability density function. Then, this heteroscedastic NLAR model is geometrically ergodic.*

Appendix B FORWARD BOOTSTRAP PREDICTION OF GENERAL NLAR MODELS

Similar to the setup in the main context, we suppose that we observe $T + p$ number of real-valued samples $\{X_{-p+1}, X_{-p+2}, \dots, X_T\}$, but these samples were generated from the below general ergodic NLAR model:

$$X_t = G(\mathbf{X}_{t-1}, \epsilon_t); \quad (6)$$

here $\{\epsilon_t\}$ is still assumed to be *i.i.d.* with mean zero, and \mathbf{X}_{t-1} represents vector $\{X_{t-1}, \dots, X_{t-p}\}$; $G(\cdot, \cdot)$ can be any continuous (possibly non-linear) function that makes the variance and mean of $\{X_t\}$ finite. The problem is how can we make multi-step ahead prediction inferences with such a complicated model. As different from the location-scale model in the main context, it may be hard to estimate F_ϵ by residuals explicitly in practice. Thus, we present forward bootstrap algorithms for two ideal cases: (a) $G(\cdot, \cdot)$ and F_ϵ are known; (b) $G(\cdot, \cdot)$ and F_ϵ are unknown, but they can be estimated consistently.

B.1 Simulation-based prediction

First, consider an idealized situation where the model $G(\mathbf{X}_{t-1}, \epsilon_t)$ and the distribution of innovations F_ϵ are known. In this case, we propose to approximate the h -step ahead prediction by simulating innovations from F_ϵ and plugging them into the NLAR model. To describe the idea, focus on the

2-step ahead prediction, and note that the distribution of the future value X_{T+2} (conditional on the observed \mathbf{X}_T) is identical to the distribution of

$$G(G(\mathbf{X}_T, \epsilon_{T+1}^*), \epsilon_{T+2}^*),$$

where ϵ_{T+1}^* and ϵ_{T+2}^* are *i.i.d.* $\sim F_\epsilon$. Going to the h -step ahead prediction, the distribution of the future value X_{T+h} (conditional on the observed \mathbf{X}_T) is identical to the distribution of the quantity

$$G(\cdots G(G(G(\mathbf{X}_T, \epsilon_{T+1}^*), \epsilon_{T+2}^*), \epsilon_{T+3}^*), \dots, \epsilon_{T+h}^*), \quad (7)$$

where $\epsilon_{T+1}^*, \dots, \epsilon_{T+h}^*$ are *i.i.d.* $\sim F_\epsilon$.

Of course, in order to obtain the L_2 or L_1 optimal point predictor, we would need to approximate the mean or median of the quantity (7). We can do this by Monte Carlo (MC) simulation; the simulation will be based on M replicates of the quantity (7); these are denoted $\{X_{T+h}^{(m)}\}_{m=1}^M$. Then, the L_2 or L_1 optimal predictor of X_{T+h} can be approximated by the mean or median, respectively, of $\{X_{T+h}^{(1)}, \dots, X_{T+h}^{(M)}\}$; the empirical distribution of the values $\{X_{T+h}^{(1)}, \dots, X_{T+h}^{(M)}\}$ can also be used to approximate the distribution of the future value X_{T+h} (conditional on the observed \mathbf{X}_T), leading to the construction of asymptotically valid PIs as $M \rightarrow \infty$.

B.2 Bootstrap-based prediction

In the more realistic scenario, both $G(\cdot, \cdot)$ and F_ϵ are unknown but we assume that they can be estimated from the data at hand; denote their estimators by $\hat{G}(\cdot, \cdot)$ and the empirical distribution of residuals by \hat{F}_ϵ , respectively, and assume they are consistent. Conducting a simulation as described in the previous subsection using $\hat{G}(\cdot, \cdot)$ and \hat{F}_ϵ in place of the unknown $G(\cdot, \cdot)$ and F_ϵ turns the MC simulation into a *bootstrap* procedure. The bootstrap version of the quantity (7) is now given by

$$\hat{G}(\cdots \hat{G}(\hat{G}(\hat{G}(\mathbf{X}_T, \hat{\epsilon}_{T+1}^*), \hat{\epsilon}_{T+2}^*), \hat{\epsilon}_{T+3}^*), \dots, \hat{\epsilon}_{T+h}^*), \quad (8)$$

where $\{\hat{\epsilon}_t^*\}_{t=T+1}^{T+h}$ are *i.i.d.* $\sim \hat{F}_\epsilon$; $\hat{G}(\cdot, \cdot)$ is an estimator to the true model. Next, we can take a similar approach as the simulation-based method previously described to approximate the L_1 or L_2 optimal predictor of X_{T+h} . Also, the asymptotically valid PI can be determined as $T \rightarrow \infty$ and $M \rightarrow \infty$. However, such a QPI will suffer from finite-sample undercoverage since the variability of estimating the model is not taken into account. We will prefer constructing a PPI instead that possesses a stronger property than the asymptotic validity. The setup of PPI for general NLAR models is similar to the procedure explained by Algorithm 2 in the main text.

Appendix C PROOFS

PROOF OF THEOREM 2.1. Denote the exactly L_2 optimal predictor of X_{T+h} by $X_{T+h}^{L_2}$. We know it is a conditional mean given observed data:

$$X_{T+h}^{L_2} = \mathbb{E}(X_{T+h} | X_T, \dots, X_{T-p+1}). \quad (9)$$

Due to the strong stationary property, we can rewrite Eq.(10) in the main text as $X_{T+h} = f(\mathbf{Y}, \boldsymbol{\epsilon}_h)$, where \mathbf{Y} and $\boldsymbol{\epsilon}_h$ represent $\{X_t\}_{t=1}^p$ and $\{\epsilon_t\}_{t=T+1}^{T+h}$, respectively. Moreover, by the causality assumption,

we get \mathbf{Y} and $\boldsymbol{\epsilon}_h$ are independent. In addition, $\{\boldsymbol{\epsilon}_h^{(i)}\}_{i=1}^M$ are also *i.i.d.*. Thus, $\{X_{T+h}^{(i)}\}_{i=1}^M$ are conditionally *i.i.d.* given \mathbf{Y} . Based on Theorem 4.2 of Majerek, Nowak, and Zieba (2005), the conditional version of the strong law of large numbers implies that:

$$\widehat{X_{T+h}^{L_2}} = \frac{1}{M} \sum_{i=1}^M X_{T+h}^{(i)} \xrightarrow{a.s.} X_{T+h}^{L_2}, \text{ assuming that } X_{T+h}^{L_2} \text{ exists.} \quad (10)$$

The existence of $X_{T+h}^{L_2}$ is guaranteed by assumptions A2 and A5 - A7. This proof can be directly extended to the NLAR model with heteroscedastic errors, since the relationship $X_{T+h} = f(\mathbf{Y}, \boldsymbol{\epsilon}_h)$ is also satisfied, so $\{X_{T+h}^{(i)}\}_{i=1}^M$ are still conditionally *i.i.d.* given \mathbf{Y} . Therefore, without changing other parts of this proof, we can show the analogous theorem for NLAR models with heteroscedastic error cases. \square

PROOF OF THEOREM 2.2. We actually want to show that under the ergodic property:

$$\widehat{X_{T+h}^{L_1}} \xrightarrow{p} X_{T+h}^{L_1}, \text{ as } M \text{ converges to infinity.} \quad (11)$$

We can write $\widehat{X_{T+h}^{L_1}}$ and $X_{T+h}^{L_1}$ as:

$$\widehat{X_{T+h}^{L_1}} = H_M(\mathbf{Y}, \mathcal{E}_h) = \text{Median}(f(\mathbf{Y}, \boldsymbol{\epsilon}_h^{(1)}), \dots, f(\mathbf{Y}, \boldsymbol{\epsilon}_h^{(M)})) ; X_{T+h}^{L_1} = H(\mathbf{Y}, \boldsymbol{\epsilon}_h) = Q_{X_{T+h}|\mathbf{X}_T}(1/2), \quad (12)$$

where \mathbf{Y} and $\boldsymbol{\epsilon}_h$ represent $\{X_t\}_{t=1}^p$ and $\{\epsilon_t\}_{t=T+1}^{T+h}$, respectively; \mathcal{E}_h represents the whole set $\{\boldsymbol{\epsilon}_h^{(i)}\}_{i=1}^M$. $Q_{X_{T+h}|\mathbf{X}_T}$ is the conditional quantile function of X_{T+h} . By assumption A8, $f(\cdot, \cdot)$ is also uniformly continuous in \mathbf{x} since it is a composition of uniformly continuous functions. Thus, for a given $\eta > 0$, there exists a constant $\delta > 0$ such that:

$$|f(\mathbf{y}_1, \boldsymbol{\epsilon}_h) - f(\mathbf{y}_2, \boldsymbol{\epsilon}_h)| < \eta, \text{ when } \|\mathbf{y}_1 - \mathbf{y}_2\| \leq \delta, \quad (13)$$

where $\|\cdot\|$ is any norm equivalent to the Euclidean norm.

Then, we can split the p -dimensional ball $B(D) = \{\mathbf{y}, \|\mathbf{y}\| \leq D\}$ into some disjoint subsets S_j , $j = 1, \dots, k$. Let $\|\mathbf{y} - \mathbf{s}_j\| < \delta$ for $\mathbf{s}_j \in S_j$ and any point $\mathbf{y} \in B(D)$ such that:

$$|f(\mathbf{y}, \boldsymbol{\epsilon}_h) - f(\mathbf{s}_j, \boldsymbol{\epsilon}_h)| < \eta, \forall \eta > 0. \quad (14)$$

Thus, $|f(\mathbf{y}, \boldsymbol{\epsilon}_h^{(i)}) - f(\mathbf{s}_j, \boldsymbol{\epsilon}_h^{(i)})| < \eta$, $\forall i \in 1, \dots, M$. It is possible to fix a small enough η to make sure that the order of $\{f(\mathbf{y}, \boldsymbol{\epsilon}_h^{(1)}), \dots, f(\mathbf{y}, \boldsymbol{\epsilon}_h^{(M)})\}$ is same with the order of $\{f(\mathbf{s}_j, \boldsymbol{\epsilon}_h^{(1)}), \dots, f(\mathbf{s}_j, \boldsymbol{\epsilon}_h^{(M)})\}$. In other words, we have:

$$\sum_{j=1}^k I_j |H_M(\mathbf{Y}, \mathcal{E}_h) - H_M(\mathbf{s}_j, \mathcal{E}_h)| < \eta, \quad (15)$$

where I_j represents the indicator function $I(\mathbf{Y} \in S_j)$, $j = 1, \dots, k$. In addition, we also have:

$$\sum_{j=1}^k I_j |H(\mathbf{Y}, \boldsymbol{\epsilon}_h) - H(\mathbf{s}_j, \boldsymbol{\epsilon}_h)| < \eta. \quad (16)$$

Therefore, define $I_0 := I(\mathbf{Y} \notin B(D))$, by combining Eqs. (15) and (16), we can get:

$$\begin{aligned}
& |H_M(\mathbf{Y}, \mathcal{E}_h) - H(\mathbf{Y}, \epsilon_h)| \\
&= \sum_{j=1}^k I_j |H_M(\mathbf{Y}, \mathcal{E}_h) - H_M(\mathbf{s}_j, \mathcal{E}_h) + H_M(\mathbf{s}_j, \mathcal{E}_h) - H(\mathbf{s}_j, \epsilon_h) + H(\mathbf{s}_j, \epsilon_h) - H(\mathbf{Y}, \epsilon_h)| + \\
& I_0 |H_M(\mathbf{Y}, \mathcal{E}_h) - H_M(\mathbf{s}_0, \mathcal{E}_h) + H_M(\mathbf{s}_0, \mathcal{E}_h) - H(\mathbf{s}_0, \epsilon_h) + H(\mathbf{s}_0, \epsilon_h) - H(\mathbf{Y}, \epsilon_h)| \\
&\leq \sum_{j=1}^k I_j (|H_M(\mathbf{Y}, \mathcal{E}_h) - H_M(\mathbf{s}_j, \mathcal{E}_h)| + |H_M(\mathbf{s}_j, \mathcal{E}_h) - H(\mathbf{s}_j, \epsilon_h)| + |H(\mathbf{s}_j, \epsilon_h) - H(\mathbf{Y}, \epsilon_h)|) + I_0 \cdot C \\
&\leq 2\eta + I_0 \cdot C + \sum_{j=1}^k I_j |H_M(\mathbf{s}_j, \mathcal{E}_h) - H(\mathbf{s}_j, \epsilon_h)|.
\end{aligned} \tag{17}$$

Comparing $H_M(\mathbf{s}_j, \mathcal{E}_h)$ and $H(\mathbf{s}_j, \epsilon_h)$, where \mathbf{s}_j is a fixed point, by applying the CLT on the sample median for the ergodic series, we can get $H_M(\mathbf{s}_j, \mathcal{E}_h)$ converges to $H(\mathbf{s}_j, \epsilon_h)$ in probability. The non-zero property of the probability density of X_{n+h} at the median is guaranteed by the everywhere positive density function of innovation. Thus, we have:

$$\mathbb{P}(|H_M(\mathbf{Y}, \mathcal{E}_h) - H(\mathbf{Y}, \epsilon_h)| \leq 2\eta + I_0 \cdot C) \rightarrow 1. \tag{18}$$

Besides, D can be arbitrarily large. Also, η can be arbitrarily small. Finally, we get:

$$\mathbb{P}(|H_M(\mathbf{Y}, \mathcal{E}_h) - H(\mathbf{Y}, \epsilon_h)| \leq \varepsilon) \rightarrow 1. \tag{19}$$

The above proof can be extended to other quantile estimators. Thus, we can build asymptotically valid QPI with any CVR. For extending such proof to NLAR models with heteroscedastic errors, we need the variance function is also uniformly continuous. Then, $f(\cdot, \cdot)$ is still uniformly continuous in \mathbf{x} . Therefore, without changing other parts of this proof, we can show the analogous theorem for NLAR models with heteroscedastic error cases. \square

PROOF OF LEMMA 2.1. To simplify the notation, we just consider the NLAR model in Eq.(11) of the main text with order 1 and homoscedastic errors. The case with higher order and heteroscedastic errors can be proved similarly.

The proof is inspired by the work of Boldin (1983). For connecting $F_\epsilon(x)$ and $\hat{F}_\epsilon(x)$, we should notice that the empirical distribution $F_T(x)$ based on $\{\epsilon_i\}_{i=1}^T$ can be a bridge, i.e., we have:

$$\begin{aligned}
\sup_x |\hat{F}_\epsilon(x) - F_\epsilon(x)| &= \sup_x |\hat{F}_\epsilon(x) - F_T(x) + F_T(x) - F_\epsilon(x)| \\
&\leq \sup_x |\hat{F}_\epsilon(x) - F_T(x)| + \sup_x |F_T(x) - F_\epsilon(x)|.
\end{aligned} \tag{20}$$

From the Glivenko–Cantelli theorem, we know $\sup_x |F_T(x) - F_\epsilon(x)|$ converges to 0 a.s.. Thus, we only need to show:

$$\sup_x |\hat{F}_\epsilon(x) - F_T(x)| \xrightarrow{P} 0. \tag{21}$$

First, we know

$$\text{if } \Delta_i(x) = \begin{cases} 1, & \epsilon_i \leq x \\ 0, & \epsilon_i > x \end{cases} \text{ then, } F_T(x) = \frac{1}{T} \sum_{i=1}^T \Delta_i(x). \tag{22}$$

Thus, we can get:

$$\widehat{F}_\epsilon(x) = \frac{1}{T} \sum_{i=1}^T \Delta_i(x + \phi(X_{i-1}, \widehat{\theta}_1) - \phi(X_{i-1}, \theta_1)), \quad (23)$$

since $\widehat{\epsilon}_i = \phi(X_{i-1}, \theta_1) - \phi(X_{i-1}, \widehat{\theta}_1) + \epsilon_i$. For handling the randomness of $\phi(X_{i-1}, \widehat{\theta}_1) - \phi(X_{i-1}, \theta_1)$ inside $\Delta_i(\cdot)$ of Eq. (23), we use nonrandom $\eta_T, T = 1, 2, 3, \dots$, to replace $\phi(X_{i-1}, \widehat{\theta}_1) - \phi(X_{i-1}, \theta_1)$. Then, we can consider the process:

$$z_T(x, \eta_T) = \widehat{F}_\epsilon(x) - F_T(x) = \frac{1}{T} \sum_{i=1}^T (\Delta_i(x + \eta_T) - \Delta_i(x)). \quad (24)$$

Indeed, we have:

$$\mathbb{P}(\sup_x |\widehat{F}_\epsilon(x) - F_T(x)| > \epsilon) \leq \mathbb{P}(\sup_x \sup_{|\eta_T| \leq T^{-\lambda}} |z_T(x, \eta_T)| > \epsilon) + \mathbb{P}(|\phi(X_{i-1}, \widehat{\theta}_1) - \phi(X_{i-1}, \theta_1)| > T^{-\lambda}). \quad (25)$$

Without loss of generality, we select an appropriate λ to make sure the second term on the right-hand side of the above inequality converges to 0 under A9 and A10.

Then, we shall show $\mathbb{P}(\sup_x \sup_{|\eta_T| \leq T^{-\lambda}} |z_T(x, \eta_T)| > \epsilon)$ also converges to 0. Since this term depends on the continuum of values of x , we can partition the real axis into $N_T \sim T^{1/2}$ parts by points:

$$-\infty = x_0 < x_1 < \dots < x_k < \dots < x_{N_T-1} < x_{N_T} = \infty, \text{ where } F_\epsilon(x_k) = kN_T^{-1}. \quad (26)$$

Hence, for x_r and x_{r+1} such that $x_r \leq x \leq x_{r+1}$, we have:

$$x_r + \eta_T \leq x + \eta_T \leq x_{r+1} + \eta_T. \quad (27)$$

In addition, since $\Delta_i(x)$ is monotonic, we obtain:

$$\begin{aligned} z_T(x, \eta_T) &\geq z_T(x_r, \eta_T) + \frac{1}{T} \sum_{i=1}^T \Delta_i(x_r) - \frac{1}{T} \sum_{i=1}^T \Delta_i(x_{r+1}); \\ z_T(x, \eta_T) &\leq z_T(x_{r+1}, \eta_T) + \frac{1}{T} \sum_{i=1}^T \Delta_i(x_{r+1}) - \frac{1}{T} \sum_{i=1}^T \Delta_i(x_r). \end{aligned} \quad (28)$$

Therefore, we have:

$$\begin{aligned} \sup_x \sup_{|\eta_T| \leq T^{-\lambda}} |z_T(x, \eta_T)| &\leq \sup_{k \leq N_T-1} \sup_{|\eta_T| \leq T^{-\lambda}} |z_T(x_{k+1}, \eta_T)| \\ &\quad + \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} |z_T(x_k, \eta_T)| \\ &\quad + \sup_{|t_1 - t_2| \leq N_T^{-1}} \frac{1}{T} \left| \sum_{i=1}^T (\Delta_i(F_\epsilon^{-1}(t_1)) - \Delta_i(F_\epsilon^{-1}(t_2))) \right|. \end{aligned} \quad (29)$$

For the last term on the r.h.s. of Eq. (29):

$$\begin{aligned}
& \sup_{|t_1-t_2| \leq N_T^{-1}} \frac{1}{T} \left| \sum_{i=1}^T (\Delta_i(F_\epsilon^{-1}(t_1)) - \Delta_i(F_\epsilon^{-1}(t_2))) \right| \\
&= \sup_{|t_1-t_2| \leq N_T^{-1}} \frac{1}{T} \left| \sum_{i=1}^T (\Delta_i(F_\epsilon^{-1}(t_1)) - t_1 + t_1 - t_2 + t_2 - \Delta_i(F_\epsilon^{-1}(t_2))) \right| \\
&\leq \sup_{t_1, s.t. |t_1-t_2| \leq N_T^{-1}} \left| \frac{1}{T} \sum_{i=1}^T \Delta_i(F_\epsilon^{-1}(t_1)) - t_1 \right| + \sup_{|t_1-t_2| \leq N_T^{-1}} |t_1 - t_2| + \sup_{t_2, s.t. |t_1-t_2| \leq N_T^{-1}} \left| t_2 - \frac{1}{T} \sum_{i=1}^T \Delta_i(F_\epsilon^{-1}(t_2)) \right|.
\end{aligned} \tag{30}$$

By the Glivenko–Cantelli theorem, it is obvious that $\sup_{|t_1-t_2| \leq N_T^{-1}} \frac{1}{T} \left| \sum_{i=1}^T (\Delta_i(F_\epsilon^{-1}(t_1)) - \Delta_i(F_\epsilon^{-1}(t_2))) \right|$ is $o_p(1)$. Next, we consider the second term of the r.h.s of Eq. (29):

$$\begin{aligned}
& \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} |z_T(x_k, \eta_T)| \\
&= \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} \left| \frac{1}{T} \sum_{i=1}^T (\Delta_i(x_k + \eta_T) - \Delta_i(x_k)) \right| \\
&= \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} \left| \frac{1}{T} \sum_{i=1}^T \Delta_i(x_k + \eta_T) - F_\epsilon(x_k + \eta_T) + F_\epsilon(x_k + \eta_T) - F_\epsilon(x_k) + F_\epsilon(x_k) - \frac{1}{T} \sum_{i=1}^T \Delta_i(x_k) \right| \\
&\leq \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} \left\{ \left| \frac{1}{T} \sum_{i=1}^T \Delta_i(x_k + \eta_T) - F_\epsilon(x_k + \eta_T) \right| + |F_\epsilon(x_k + \eta_T) - F_\epsilon(x_k)| + \left| F_\epsilon(x_k) - \frac{1}{T} \sum_{i=1}^T \Delta_i(x_k) \right| \right\}.
\end{aligned} \tag{31}$$

Applying the Glivenko–Cantelli theorem again, we can find the first and third term in the r.h.s. of Eq. (31) converges to 0 a.s.. For the middle term:

$$\begin{aligned}
& \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} |F_\epsilon(x_k + \eta_T) - F_\epsilon(x_k)| \\
&= \sup_{k \leq N_T} \sup_{|\eta_T| \leq T^{-\lambda}} |F_\epsilon(x_k) + F'_\epsilon(\tilde{x})\eta_T - F_\epsilon(x_k)| \\
&\leq \sup_x \sup_{|\eta_T| \leq T^{-\lambda}} |F'_\epsilon(\tilde{x})\eta_T| \rightarrow 0, \text{ Under A11.}
\end{aligned} \tag{32}$$

We can do a similar analysis for the first term of the r.h.s of Eq. (29). Combining all parts together, we prove Lemma 2.1. The proof with heteroscedastic errors can be written similarly. \square

PROOF OF THEOREM 2.3. For simplifying the notations, we consider the NLAR model with order 1 and homoscedastic errors. For NLAR models with higher order and heteroscedastic errors, the proof can be written similarly. We show the proof of $h = 2$ as an example. The proof of higher steps or one-step prediction can be written similarly.

First, by the tower property, we can show that $F_{X_{T+2}|X_T}(x)$ is equivalent to:

$$\begin{aligned}
F_{X_{T+2}|X_T}(x) &= \mathbb{P}(X_{T+2} \leq x | X_T) \\
&= \mathbb{P}(\phi(X_{T+1}, \theta_1) + \epsilon_{T+2} \leq x | X_T) \\
&= \mathbb{P}\left(\epsilon_{T+2} \leq x - \phi(\phi(X_T, \theta_1) + \epsilon_{T+1}, \theta_1) \middle| X_T\right) \\
&= \mathbb{E}\left[\mathbb{P}\left(\epsilon_{T+2} \leq x - \phi(\phi(X_T, \theta_1) + \epsilon_{T+1}, \theta_1) \middle| \epsilon_{T+1}, X_T\right) \middle| X_T\right] \\
&= \mathbb{E}\left[F_\epsilon(x - \phi(\phi(X_T, \theta_1) + \epsilon_{T+1}, \theta_1)) \middle| X_T\right] \\
&= \mathbb{E}\left[F_\epsilon(\mathcal{L}(x, X_T, \epsilon_{T+1})) \middle| X_T\right];
\end{aligned} \tag{33}$$

we use $\mathcal{L}(x, X_T, \epsilon_{T+1})$ to represent $x - \phi(\phi(X_T, \theta_1) + \epsilon_{T+1}, \theta_1)$ to simplify notations. Similarly, we can analyze $F_{X_{T+2}^*|X_T, \dots, X_0}(x)$, it has below equivalent expressions:

$$\begin{aligned}
F_{X_{T+2}^*|X_T, \dots, X_0}(x) &= \mathbb{P}(X_{T+2}^* \leq x | X_T, \dots, X_0) \\
&= \mathbb{E}\left[\mathbb{P}\left(\hat{\epsilon}_{T+2}^* \leq \widehat{\mathcal{L}}(x, X_T, \hat{\epsilon}_{T+1}^*) \middle| \hat{\epsilon}_{T+1}^*, X_T, \dots, X_0\right) \middle| X_T, \dots, X_0\right] \\
&= \mathbb{E}^*\left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, X_T, \hat{\epsilon}_{T+1}^*))\right],
\end{aligned} \tag{34}$$

where $\widehat{\mathcal{L}}(x, X_T, \hat{\epsilon}_{T+1}^*)$ represents $x - \phi(\phi(X_T, \hat{\theta}_1) + \hat{\epsilon}_{T+1}^*, \hat{\theta}_1)$ and $\mathbb{E}^*(\cdot)$ represents the expectation in the bootstrap world, i.e., $\mathbb{E}(\cdot | X_T, \dots, X_0)$. Thus, we hope to show:

$$\sup_{|x| \leq c_T} \left| \mathbb{E}^*\left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, X_T, \hat{\epsilon}_{T+1}^*))\right] - \mathbb{E}\left[F_\epsilon(\mathcal{L}(x, X_T, \epsilon_{T+1})) \middle| X_T\right] \right| \xrightarrow{p} 0. \tag{35}$$

From here, we first bound the region of X_t by Lemma 1 of Franke, Neumann, and Stockis (2004) under A1–A7:

$$\mathbb{P}(|X_T| > \gamma_T) \rightarrow 0, \tag{36}$$

where $\{\gamma_T\}$ is a sequence of sets, such that $\gamma_1 \subseteq \dots \subseteq \gamma_T \subseteq \gamma_{T+1} \subseteq \dots$ with the form $\gamma_T = [-T^c, T^c]$; c is some appropriate constant. For deriving such a result for a time series model with heteroscedastic errors, we need the additional assumption of variance function in A1; then the proof is referred to Lemma 1 of Franke, Neumann, and Stockis (2004). In addition, we have a relationship:

$$\begin{aligned}
&\mathbb{P}\left(\sup_{|x| \leq c_T} \left| \mathbb{E}^*\left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, X_T, \hat{\epsilon}_{T+1}^*))\right] - \mathbb{E}\left[F_\epsilon(\mathcal{L}(x, X_T, \epsilon_{T+1})) \middle| X_T\right] \right| > \varepsilon\right) \\
&\leq \mathbb{P}(|X_T| > \gamma_T) + \mathbb{P}\left((|X_T| \leq \gamma_T) \cap \left(\sup_{|x| \leq c_T} \left| \mathbb{E}^*\left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, X_T, \hat{\epsilon}_{T+1}^*))\right] - \mathbb{E}\left[F_\epsilon(\mathcal{L}(x, X_T, \epsilon_{T+1})) \middle| X_T\right] \right| > \varepsilon\right)\right).
\end{aligned} \tag{37}$$

Thus, to verify Eq. (35), we just need to show that the second term of the r.h.s. of Eq. (37) converges to 0. We can take the sequence c_T and γ_T to be the same sequence which converges to infinity slowly enough. Then, it is enough for us to analyze the asymptotic probability of the below expression:

$$\sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^*\left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*))\right] - \mathbb{E}\left[F_\epsilon(\mathcal{L}(x, y, \epsilon_{T+1}))\right] \right| > \varepsilon. \tag{38}$$

Decompose the l.h.s. of Eq. (38) as:

$$\begin{aligned}
& \sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^* \left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] - \mathbb{E} [F_\epsilon(\mathcal{L}(x, y, \epsilon_{T+1}))] \right| \\
& \leq \sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^* \left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] - \mathbb{E}^* \left[F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] \right| \\
& \quad + \sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^* \left[F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] - \mathbb{E} [F_\epsilon(\mathcal{L}(x, y, \epsilon_{T+1}))] \right|.
\end{aligned} \tag{39}$$

Then, we analyze two terms on the r.h.s. of Eq. (39) separately. For the first term, we have:

$$\begin{aligned}
& \sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^* \left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] - \mathbb{E}^* \left[F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] \right| \\
& \leq \sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^* \left[\widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) - F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] \right| \\
& \leq \sup_{|x|, |y| \leq c_T, z} \left| \widehat{F}_\epsilon(\widehat{\mathcal{L}}(x, y, z)) - F_\epsilon(\widehat{\mathcal{L}}(x, y, z)) \right| \xrightarrow{p} 0, \text{ under Lemma 4.1.}
\end{aligned} \tag{40}$$

For the second term on the r.h.s. of Eq. (39), we have:

$$\begin{aligned}
& \sup_{|x|, |y| \leq c_T} \left| \mathbb{E}^* \left[F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_{T+1}^*)) \right] - \mathbb{E} [F_\epsilon(\mathcal{L}(x, y, \epsilon_{T+1}))] \right| \\
& \leq \sup_{|x|, |y| \leq c_T} \left| \frac{1}{T} \sum_{i=1}^T F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i)) - \frac{1}{T} \sum_{i=1}^T F_\epsilon(\mathcal{L}(x, y, \epsilon_i)) \right| \\
& \quad + \sup_{|x|, |y| \leq c_T} \left| \frac{1}{T} \sum_{i=1}^T F_\epsilon(\mathcal{L}(x, y, \epsilon_i)) - \mathbb{E} [F_\epsilon(\mathcal{L}(x, y, \epsilon_{T+1}))] \right|,
\end{aligned} \tag{41}$$

where $\{\epsilon_j\}_{j=1}^T$ are taken as $X_i - \phi(X_{i-1}, \theta_1)$ for $i = 1, \dots, T$. We can show:

$$\begin{aligned}
& \mathbb{P} \left(\max_{i=1, \dots, T} |\epsilon_i - \hat{\epsilon}_i| > \varepsilon \right) \\
& = \mathbb{P} \left(\max_{i=1, \dots, T} \left| X_i - \phi(X_{i-1}, \theta_1) - X_i + \phi(X_{i-1}, \hat{\theta}_1) \right| > \varepsilon \right) \\
& \leq \mathbb{P} \left(\max_{i=1, \dots, T} |X_{i-1}| > c_T \right) \\
& \quad + \mathbb{P} \left(\left(\max_{i=1, \dots, T} |X_{i-1}| < c_T \right) \cap \left(\max_{i=1, \dots, T} \left| \phi(X_{i-1}, \hat{\theta}_1) - \phi(X_{i-1}, \theta_1) \right| > \varepsilon \right) \right) \\
& \leq o(1) + \mathbb{P} \left(\sup_{|x| \leq c_T} \left| \phi(x, \hat{\theta}_1) - \phi(x, \theta_1) \right| > \varepsilon \right) \\
& \rightarrow 0, \text{ under A9 and A10.}
\end{aligned} \tag{42}$$

We further consider two terms on the r.h.s. of Eq. (41) separately. For the first term, by Taylor

expansion, we have:

$$\begin{aligned}
& \sup_{|x|, |y| \leq c_T} \left| \frac{1}{T} \sum_{i=1}^T F_\epsilon(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i)) - \frac{1}{T} \sum_{i=1}^T F_\epsilon(\mathcal{L}(x, y, \epsilon_i)) \right| \\
&= \sup_{|x|, |y| \leq c_T} \left| \frac{1}{T} \sum_{i=1}^T \left(F_\epsilon(\mathcal{L}(x, y, \epsilon_i)) + f_\epsilon(o_i)(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i) - \mathcal{L}(x, y, \epsilon_i)) \right) - \frac{1}{T} \sum_{i=1}^T F_\epsilon(\mathcal{L}(x, y, \epsilon_i)) \right| \\
&= \sup_{|x|, |y| \leq c_T} \left| \frac{1}{T} \sum_{i=1}^T f_\epsilon(o_i)(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i) - \mathcal{L}(x, y, \epsilon_i)) \right| \\
&\leq \sup_{|x|, |y| \leq c_T} \frac{1}{T} \sum_{i=1}^T \left| f_\epsilon(o_i)(\widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i) - \mathcal{L}(x, y, \epsilon_i)) \right| \tag{43} \\
&\leq \sup_{|x|, |y| \leq c_T} \sup_z |f_\epsilon(z)| \cdot \frac{1}{T} \sum_{i=1}^T \left| \widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i) - \mathcal{L}(x, y, \epsilon_i) \right| \\
&\leq \sup_{|x|, |y| \leq c_T} C \cdot \frac{1}{T} \sum_{i=1}^T \left| \widehat{\mathcal{L}}(x, y, \hat{\epsilon}_i) - \mathcal{L}(x, y, \epsilon_i) \right| \text{ (under A11)} \\
&\leq \sup_{|x|, |y| \leq c_T, j \in \{1, \dots, T\}} C \cdot \left| \widehat{\mathcal{L}}(x, y, \hat{\epsilon}_j) - \mathcal{L}(x, y, \epsilon_j) \right|.
\end{aligned}$$

From Eq. (42) and A9–A12, we have Eq. (43) converges to 0 in probability. For the second term on the r.h.s. of Eq. (41), by the uniform law of large numbers, we have:

$$\sup_{|x|, |y| \leq c_T} \left| \frac{1}{T} \sum_{i=1}^T F_\epsilon(\mathcal{L}(x, y, \epsilon_i)) - \mathbb{E}[F_\epsilon(\mathcal{L}(x, y, \epsilon_{T+1}))] \right| \xrightarrow{P} 0. \tag{44}$$

Combine all pieces, Eq. (38) converges to 0 in probability, which implies Theorem 2.3. \square

PROOF OF LEMMA 3.1. Under A4–A7, the time series $\{X_t\}$ is ergodic. Besides, under A2 and A10, we can show $L(\vartheta)$ is uniformly finite for $\vartheta \in \Theta_1$. With A13, by the uniform law of larger numbers for the ergodic series, see Theorem 6 of Kirch and Kamgaing (2012) for a reference, we have:

$$\sup_{\vartheta \in \Theta_1} |L_T(\vartheta) - L(\vartheta)| \xrightarrow{P} 0. \tag{45}$$

Under A13 and A14, we can easily show that:

$$\inf_{|\vartheta - \theta_1| > \epsilon} L(\vartheta) > L(\theta_1), \text{ for } \forall \epsilon, \tag{46}$$

Eq. (46) implies that given $\forall \epsilon > 0$, $\exists C > 0$ such that $|\vartheta - \theta_1| > \epsilon \Rightarrow L(\vartheta) - L(\theta_1) \geq C > 0$, thus we have:

$$\begin{aligned}
\mathbb{P}(|\hat{\theta}_1 - \theta_1| > \epsilon) &\leq \mathbb{P}(L(\hat{\theta}_1) - L(\theta_1) \geq C) \\
&= \mathbb{P}(L(\hat{\theta}_1) - L_T(\hat{\theta}_1) + L_T(\hat{\theta}_1) - L(\theta_1) \geq C) \\
&\leq \mathbb{P}(L(\hat{\theta}_1) - L_T(\hat{\theta}_1) + L_T(\theta_1) - L(\theta_1) \geq C) \\
&\leq \mathbb{P}(2 \sup_{\vartheta \in \Theta_1} |L_T(\vartheta) - L(\vartheta)| > C) \rightarrow 0.
\end{aligned} \tag{47}$$

The last inequality of Eq. (47) is guaranteed by Eq. (45). \square

PROOF OF LEMMA 3.2. The proof of the consistency of $\hat{\theta}_1$ to θ_1 is the same as the proof of Lemma 4.2. Similar to the proof of Lemma 4.2, by the ergodic property of the series $\{X_t\}$, we know:

$$\sup_{\vartheta \in \Theta_2} \left| \frac{1}{T} \sum_{t=1}^T \left(\frac{X_t - \phi(X_{t-1}, \theta_1)}{h(X_{t-1}, \vartheta)} \right)^2 - \mathbb{E} \left(\frac{X_t - \phi(X_{t-1}, \theta_1)}{h(X_{t-1}, \vartheta)} \right)^2 \right| \xrightarrow{p} 0. \quad (48)$$

Actually, after the first step estimation, we have got the $\hat{\theta}_1$ which is consistent to the true parameter θ_1 . Thus, we need to find the below convergence relationship:

$$\sup_{\vartheta \in \Theta_2} \left| \frac{1}{T} \sum_{t=1}^T \left(\frac{X_t - \phi(X_{t-1}, \hat{\theta}_1)}{h(X_{t-1}, \vartheta)} \right)^2 - \mathbb{E} \left(\frac{X_t - \phi(X_{t-1}, \theta_1)}{h(X_{t-1}, \vartheta)} \right)^2 \right| \xrightarrow{p} 0. \quad (49)$$

Since $\hat{\theta}_1$ converges to θ_1 in probability, it is easily to find:

$$\sup_{\vartheta \in \Theta_2} \left| \frac{1}{T} \sum_{t=1}^T \left(\frac{X_t - \phi(X_{t-1}, \hat{\theta}_1)}{h(X_{t-1}, \vartheta)} \right)^2 - \frac{1}{T} \sum_{t=1}^T \left(\frac{X_t - \phi(X_{t-1}, \theta_1)}{h(X_{t-1}, \vartheta)} \right)^2 \right| \xrightarrow{p} 0, \quad (50)$$

which implies Eq. (49) for a compact set Θ_2 in conjunction with Eq. (48). Then, since $|\cdot - a|$ is a uniformly continuous function for a constant a , by applying the uniform continuous mapping theorem which is the Theorem 1 of Kasy (2019) on Eq. (49), we can get:

$$\sup_{\vartheta \in \Theta_2} |K_T(\vartheta, \hat{\theta}_1) - K(\vartheta, \theta_1)| \xrightarrow{p} 0, \quad (51)$$

where $K(\vartheta, \theta_1)$ and $K_T(\vartheta, \hat{\theta}_1)$ are:

$$K(\vartheta, \theta_1) = \left| \mathbb{E} \left(\frac{X_t - \phi(X_{t-1}, \theta_1)}{h(X_{t-1}, \vartheta)} \right)^2 - 1 \right|; \quad K_T(\vartheta, \hat{\theta}_1) = \left| \frac{1}{T} \sum_{t=1}^T \left(\frac{X_t - \phi(X_{t-1}, \hat{\theta}_1)}{h(X_{t-1}, \vartheta)} \right)^2 - 1 \right|. \quad (52)$$

Then, we can repeat the procedure in the proof of Lemma 3.1 to show:

$$\mathbb{P}(|\hat{\theta}_2 - \theta_2| > \varepsilon) \rightarrow 0, \quad \forall \varepsilon > 0. \quad (53)$$

□

PROOF OF THEOREM 3.1. Throughout this proof, we focus on a sequence of sets $\Omega_{T+p} \subseteq \mathbb{R}^{T+p}$, such that $\mathbb{P}((X_{-p+1}, \dots, X_T) \notin \Omega_{T+p}) = o(1)$. We first explain the truth that the bootstrap series is ergodic for $(X_{-p+1}, \dots, X_T) \in \Omega_{T+p}$. Particularly, we want to check analogous A4-A7 in the bootstrap world. With the consistency property of parameter estimators, i.e., A9 and the continuous density function of residuals after the convolutional manipulation, we can easily find A4-A7 are satisfied in the bootstrap world; see Theorem 2 of Franke, Neumann, and Stockis (2004) for more discussions.

To show the closeness of the stationary distribution of bootstrap and real series. Based on the proof of Theorem 3 of Franke, Kreiss, et al. (2002) and Theorem 3 of Franke, Neumann, and Stockis (2004), we can get:

$$\sup_B |\Pi(B) - \Pi^*(B)| = o(1), \quad (54)$$

which holds for all measurable sets B , where $\Pi(B)$ and $\Pi^*(B)$ represent stationary distribution for real series and bootstrap series, respectively. This implies, in conjunction with Eq.(28) in the main text and the condition that $\mathbb{P}(X_t \notin \gamma_T) = o(1)$, we can get:

$$\begin{aligned} & \int_{\mathbb{R}^2} (x_1 - \phi(x_0, \Xi_\theta))^2 \hat{f}_\epsilon(x_1 - \phi(x_0, \hat{\theta})) \pi^*(x_0) dx_1 dx_0 \\ &= \int_{\gamma_T^2} (x_1 - \phi(x_0, \Xi_\theta))^2 \hat{f}_\epsilon(x_1 - \phi(x_0, \hat{\theta})) \pi^*(x_0) dx_1 dx_0 + o(1) \\ &\rightarrow \mathbb{E}[X_1 - \phi(X_0, \Xi_\theta)]^2, \end{aligned} \quad (55)$$

which implies Theorem 3.1 directly. □

PROOF OF THEOREM 3.2. To simplify notation, we consider the case $p = q = 1$. Higher order cases can be handled similarly. We first build the estimation inference for $\hat{\theta}_1$. Starting from A19, we have:

$$op(T^{-1/2}) = \nabla L_T(\hat{\theta}_1) = \nabla L_T(\theta_1) + \nabla^2 L_T(\tilde{\theta}_1)(\hat{\theta}_1 - \theta_1), \quad (56)$$

where $\tilde{\theta}_1$ is between $\hat{\theta}_1$ and θ_1 ; hence $\tilde{\theta}_1$ also converges to θ_1 in probability. First, we consider $\nabla^2 L_T(\tilde{\theta}_1)$, which has a form as below:

$$\nabla^2 L_T(\tilde{\theta}_1) = \frac{1}{T} \nabla^2 \sum_{t=1}^T \left(X_t - \phi(X_{t-1}, \tilde{\theta}_1) \right)^2 = \frac{1}{T} \sum_{t=1}^T \nabla^2 q_t(\tilde{\theta}_1). \quad (57)$$

Under A18, since it is easy to check $\mathbb{E} \sup_{\vartheta \in \Theta_1^0} \|\nabla^2 q_1(\vartheta)\| < \infty$, by combining the uniform law of larger numbers for the ergodic series, dominated convergence theorem, the consistency between $\tilde{\theta}_1$ and θ_1 and the continuity of $L_T(\cdot)$ w.r.t. θ_1 , we can get:

$$\nabla^2 L_T(\tilde{\theta}_1) \xrightarrow{p} B_1, \quad (58)$$

where $B_1 = 2 \cdot \mathbb{E} \left(\nabla \phi(X_0, \theta_1) (\nabla \phi(X_0, \theta_1))^\top \right)$. Thus, we can multiply both side of Eq. (56) by \sqrt{T} to get:

$$op(1) = \sqrt{T} \nabla L_T(\theta_1) + (B_1 + op(1)) \sqrt{T} (\hat{\theta}_1 - \theta_1), \quad (59)$$

which further implies that:

$$\sqrt{T} (\hat{\theta}_1 - \theta_1) = -B_1^{-1} \sqrt{T} \nabla L_T(\theta_1) + op(1), \quad (60)$$

where $\sqrt{T} \nabla L_T(\theta_1)$ has a concrete form as below:

$$-\sqrt{T} \nabla L_T(\theta_1) = \frac{2}{\sqrt{T}} \sum_{t=1}^T (X_t - \phi(X_{t-1}, \theta_1)) \nabla \phi(X_{t-1}, \theta_1). \quad (61)$$

We need the CLT for strongly mixing processes to show the asymptotic normality of Eq. (61). Based on Theorem 1.7 of Bosq (2012) for univariate case, we can show the normality with the Cramér–Wold theorem. Also, Kirch and Kamgaing (2012) applied the strong invariance principle from Kuelbs and Philipp (1980) to show the asymptotic distribution of multivariate cases directly. Since we need to

analyze the asymptotic distribution of parameters in the bootstrap world, we take the first approach which is more clear. We consider:

$$-\sqrt{T}\mathbf{a}^\top \nabla L_T(\theta_1) = \frac{2}{\sqrt{T}} \sum_{t=1}^T (X_t - \phi(X_{t-1}, \theta_1)) \mathbf{a}^\top \nabla \phi(X_{t-1}, \theta_1), \quad (62)$$

where \mathbf{a}^\top is any real vector that has the same dimension as θ_1 . Since we assume we can correctly specify the model, we have:

$$-\sqrt{T}\mathbf{a}^\top \nabla L_T(\theta_1) \xrightarrow{d} N(0, \tau_1^2), \quad (63)$$

where τ_1^2 has the form as follows:

$$\begin{aligned} \tau_1^2 &= \sum_{i=-\infty}^{\infty} \text{Cov} \left(2 \cdot (X_1 - \phi(X_0, \theta_1)) \mathbf{a}^\top \nabla \phi(X_0, \theta_1), 2 \cdot (X_{i+1} - \phi(X_i, \theta_1)) \mathbf{a}^\top \nabla \phi(X_i, \theta_1) \right) \\ &= 4 \cdot \mathbb{E} \left((\mathbf{a}^\top \nabla \phi(X_0, \theta_1)) (X_1 - \phi(X_0, \theta_1))^2 (\nabla \phi(X_0, \theta_1)^\top \mathbf{a}) \right) \\ &\quad + 2 \sum_{i=1}^{\infty} \text{Cov} \left(2 \cdot (X_1 - \phi(X_0, \theta_1)) \mathbf{a}^\top \nabla \phi(X_0, \theta_1), 2 \cdot (X_{i+1} - \phi(X_i, \theta_1)) \mathbf{a}^\top \nabla \phi(X_i, \theta_1) \right) \\ &= 4 \cdot \mathbb{E}(\mathbf{a}^\top \sigma(X_0, \theta_2) \nabla \phi(X_0, \theta_1) \nabla \phi(X_0, \theta_1)^\top \sigma(X_0, \theta_2) \mathbf{a}) \\ &\quad + 2 \sum_{i=1}^{\infty} \text{Cov} \left(2 \cdot \epsilon_1 \mathbf{a}^\top \nabla \phi(X_0, \theta_1), 2 \cdot \epsilon_{i+1} \mathbf{a}^\top \nabla \phi(X_i, \theta_1) \right) \\ &= 4 \cdot \mathbb{E}(\mathbf{a}^\top \sigma(X_0, \theta_2) R_1 \sigma(X_0, \theta_2) \mathbf{a}), \end{aligned} \quad (64)$$

where $R_1 = \nabla \phi(X_0, \theta_1) \nabla \phi(X_0, \theta_1)^\top$. Thus, applying Cramér–Wold theorem, we have:

$$-\sqrt{T} \nabla L_T(\theta_1) \xrightarrow{d} N(0, \Omega_1), \quad (65)$$

where Ω_1 is $4 \cdot \mathbb{E}(\sigma(X_0, \theta_2) R_1 \sigma(X_0, \theta_2))$. Thus, we can conclude from [Eq. \(60\)](#) that:

$$\sqrt{T}(\hat{\theta}_1 - \theta_1) \xrightarrow{d} N(0, B_1^{-1} \Omega_1 B_1^{-1}). \quad (66)$$

Compared to the result in Kirch and Kamgaing (2012), they got another form of Ω_1 :

$$\Omega_1 = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\left(\nabla \sum_{t=1}^T (X_t - \phi(X_{t-1}, \theta_1))^2 \right) \left(\nabla \sum_{t=1}^T (X_t - \phi(X_{t-1}, \theta_1))^2 \right)^\top \right]. \quad (67)$$

With a correctly specified model, it is equivalent to our form.

We can also analyze the distribution of parameter estimation $\hat{\theta}_2$. By A19, we first have:

$$op(T^{-1/2}) = \nabla K_T(\theta_2, \hat{\theta}_1) + \nabla^2 K_T(\tilde{\theta}_2, \hat{\theta}_1)(\hat{\theta}_2 - \theta_2). \quad (68)$$

For simplifying the notation, we write:

$$\begin{aligned} K_T(\tilde{\theta}_2, \hat{\theta}_1) &= \left(\frac{1}{T} \sum_{t=1}^T \left(\frac{X_t - \phi(X_{t-1}, \hat{\theta}_1)}{h(X_{t-1}, \tilde{\theta}_2)} \right)^2 - 1 \right)^2 \\ &= \left(\frac{1}{T} \sum_{t=1}^T g(X_t, X_{t-1}, \tilde{\theta}_2, \hat{\theta}_1) - 1 \right)^2. \end{aligned} \quad (69)$$

We can find $\nabla K_T(\theta_2, \theta_1)$ and $\nabla^2 K_T(\theta_2, \theta_1)$ w.r.t θ_2 have a following forms respectively:

$$\begin{aligned}\nabla K_T(\theta_2, \theta_1) &= 2 \cdot \left(\frac{1}{T} \sum_{t=1}^T g(X_t, X_{t-1}, \theta_2, \theta_1) - 1 \right) \cdot \left(\frac{1}{T} \sum_{t=1}^T \nabla g(X_t, X_{t-1}, \theta_2, \theta_1) \right) \\ \nabla^2 K_T(\theta_2, \theta_1) &= 2 \cdot \left(\frac{1}{T} \sum_{t=1}^T \nabla g(X_t, X_{t-1}, \theta_2, \theta_1) \right) \cdot \left(\frac{1}{T} \sum_{t=1}^T \nabla g(X_t, X_{t-1}, \theta_2, \theta_1) \right)^\top \\ &\quad + 2 \cdot \left(\frac{1}{T} \sum_{t=1}^T g(X_t, X_{t-1}, \theta_2, \theta_1) - 1 \right) \cdot \left(\frac{1}{T} \sum_{t=1}^T \nabla^2 g(X_t, X_{t-1}, \theta_2, \theta_1) \right).\end{aligned}\tag{70}$$

Similarly to analyze $\nabla^2 L_T(\tilde{\theta}_1)$, under the consistence relationship between $\tilde{\theta}_2$ and θ_2 , $\hat{\theta}_1$ and θ_1 , we can get:

$$\nabla^2 K_T(\tilde{\theta}_1, \hat{\theta}_1) \xrightarrow{P} B_2.\tag{71}$$

where $B_2 = 2 \cdot \mathbb{E}(\nabla g(X_1, X_0, \theta_2, \theta_1)) \cdot \mathbb{E}(\nabla g(X_1, X_0, \theta_2, \theta_1)^\top)$. Looking back to [Eq. \(68\)](#), it is left to analyze $\sqrt{T} \nabla K_T(\theta_2, \hat{\theta}_1)$. Through Slutsky's theorem, it is asymptotically equivalent to analyze $\sqrt{T} \frac{2}{T} \sum_{t=1}^T [g(X_t, X_{t-1}, \theta_2, \theta_1) - 1] \cdot B_3$; $B_3 = \mathbb{E}(\nabla g(X_1, X_0, \theta_2, \theta_1))$. Apply the same technique as we analyzed the distribution of $\hat{\theta}_1$, we can get:

$$\sqrt{T} \mathbf{a}^\top \nabla K_T(\theta_2, \hat{\theta}_1) \xrightarrow{d} N(0, \tau_2^2),\tag{72}$$

where $\tau_2^2 = 4 \cdot \mathbb{E}(\mathbf{a}^\top B_3 R_2 B_3^\top \mathbf{a})$; $R_2 = (g(X_1, X_0, \theta_2, \theta_1) - 1)^2$. Thus, we have:

$$\sqrt{T} \nabla K_T(\theta_2, \hat{\theta}_1) \xrightarrow{d} N(0, \Omega_2),\tag{73}$$

where $\Omega_2 = 4 \cdot \mathbb{E}(B_3 R_2 B_3^\top)$. This further implies that:

$$\sqrt{T}(\hat{\theta}_2 - \theta_2) \xrightarrow{d} N(0, B_2^{-1} \Omega_2 B_2^{-1}).\tag{74}$$

□

References

- An H, Huang F. 1996. The geometrical ergodicity of nonlinear autoregressive models. *Statistica Sinica* **6**: 943–956.
- Boldin M. 1983. Estimation of the distribution of noise in an autoregression scheme. *Theory of Probability & Its Applications* **27**: 866–871.
- Bosq D. 2012. *Nonparametric statistics for stochastic processes: estimation and prediction*. Vol. 110. Springer Science & Business Media.
- Doob JL. 1953. *Stochastic processes*. Wiley New York.
- Franke J, Kreiss JP, et al. 2002. Properties of the nonparametric autoregressive bootstrap. *Journal of Time Series Analysis* **23**: 555–585.
- Franke J, Neumann MH, Stockis JP. 2004. Bootstrapping nonparametric estimators of the volatility function. *Journal of Econometrics* **118**: 189–218.
- Kasy M. 2019. Uniformity and the delta method. *Journal of Econometric Methods* **8**: 1–19.
- Kirch C, Kamgaing JT. 2012. Testing for parameter stability in nonlinear autoregressive models. *Journal of Time Series Analysis* **33**: 365–385.
- Kuelbs J, Philipp W. 1980. Almost sure invariance principles for partial sums of mixing B-valued random variables. *The Annals of Probability*: 1003–1036.
- Majerek D, Nowak W, Zieba W. 2005. Conditional strong law of large number. *Int. J. Pure Appl. Math* **20**: 143–156.
- Min C, Hongzhi A. 1999. The probabilistic properties of the nonlinear autoregressive model with conditional heteroskedasticity. *Acta Mathematicae Applicatae Sinica* **15**: 9–17.
- Tjøstheim D. 1990. Non-linear time series and Markov chains. *Advances in Applied Probability* **22**: 587–611.
- Tong H. 1990. *Non-linear Time Series: A Dynamical System Approach*. Oxford university press.
- Tweedie RL. 1975. Sufficient conditions for ergodicity and recurrence of Markov chains on a general state space. *Stochastic Processes and Their Applications* **3**: 385–403.