

# Generalized threshold latent variable model

Yuanbo Li

*University of International Business and Economics,*  
*e-mail: [jacoblyb@gmail.com](mailto:jacoblyb@gmail.com)*

Xunze Zheng and Chun Yip Yau

*Chinese University of Hong Kong,*  
*e-mail: [s1155002000@sta.cuhk.edu.hk](mailto:s1155002000@sta.cuhk.edu.hk); [cyyau@sta.cuhk.edu.hk](mailto:cyyau@sta.cuhk.edu.hk)*

**Abstract:** This article proposes a generalized threshold latent variable model for flexible threshold modeling of time series. The proposed model encompasses several existing models, and allows a discrete valued threshold variable. Sufficient conditions for stationarity and ergodicity are investigated. The minimum description length principle is applied to formulate a criterion function for parameter estimation and model selection. A computationally efficient procedure for optimizing the criterion function is developed based on a genetic algorithm. Consistency and weak convergence of the parameter estimates are established. Moreover, simulation studies and an application for initial public offering data are presented to illustrate the proposed methodology.

**MSC 2010 subject classifications:** 62M10.

**Keywords and phrases:** Compound Poisson process, ergodicity, genetic algorithm, minimum description length principle, multiple-threshold, piecewise modeling.

Received April 2018.

## Contents

1	Introduction . . . . .	2044
2	Generalized threshold latent variable model . . . . .	2045
3	Stationarity and ergodicity . . . . .	2047
4	Estimation and model selection criterion . . . . .	2049
5	Asymptotic inferences of GTLVM . . . . .	2053
5.1	Assumptions for asymptotic inferences . . . . .	2053
5.2	Asymptotic theorems . . . . .	2054
6	Simulation . . . . .	2056
6.1	Example 1 . . . . .	2056
6.2	Example 2 . . . . .	2058
7	Application . . . . .	2060
	Acknowledgment . . . . .	2062
	Appendix . . . . .	2063
	References . . . . .	2088

## 1. Introduction

The threshold autoregressive (TAR) model, proposed by [58], has enormous popularity in a wide range of applications. It allows the modeling of diverse behaviors under different regimes, which provides flexible descriptions of many real-world scenarios. The backgrounds, theory, applications and extensions of TAR models can be found in the excellent surveys of [59, 61].

Recently, generalizations to the TAR model have been considered by introducing threshold structures to non-linear models instead of autoregressive models. One important direction is on generalized linear models with thresholds, such as the generalized threshold mixed model (GTMM) ([53]) and the generalized threshold stochastic regression model (GTSRM) ([54]). In addition, threshold modeling has been extended to heteroscedasticity of time series, for example, the double-threshold autoregressive moving average conditional heteroskedastic (DTARMACH) model ([38]), the threshold stochastic volatility (THSV) model ([56]), the multiple-threshold double autoregressive (MTDAR) model ([33]), and the threshold double autoregressive model (TDAR) ([34]). Furthermore, threshold models with more elastic regime switching mechanisms are considered in the endogenous delay threshold model (EDTAR) of [27] and the hysteretic autoregressive (HAR) model of [36].

For the asymptotic theory of TAR models, the strong consistency and asymptotic distributions of the parameter estimates are studied by [10] and [13]. On the other hand, tests for a linear series against its threshold extension are considered in [12], [63] and [35]. For many sophisticated threshold models, asymptotic theories are developed by assuming stationarity and ergodicity of the process (for example, [54]). However, conditions for strict stationarity and ergodicity are investigated only for the self-excited threshold autoregressive (SETAR) model ([11]), the TAR model with order  $p$  ([3]), the DTARMACH model in [38], the HAR model in [36], the MTDAR model in [33], and the TDAR model in [34]. Corresponding results for the generalized linear-type threshold models remain unexplored.

More importantly, despite the well developed theoretical background of estimation theory, the estimation procedure of threshold models demands a high computational cost. Due to the irregular nature of the threshold parameters, locating the global minimum of the likelihood requires a multi-parameter grid search over all possible values of the threshold parameters, which is typically computationally infeasible; see [32] and [65]. Consequently, many threshold models are implemented assuming one threshold a priori; see [56], [53], and [54]. When the number of thresholds and the parametric model of each regime are unknown, no practical estimation method appears available except for the simplest TAR model; see [66], [15] and [16].

To tackle the above problems and further extend the flexibility, in this article we first propose the generalized threshold latent variable model (GTLVM) which covers most of the aforementioned models as special cases. In particular, the threshold variables may be continuous or discrete valued. As far as we know, asymptotic theory for threshold models with discrete-valued thresholds has not

been investigated in the literature. Second, we establish sufficient conditions for the stationarity and ergodicity of the proposed model. Third, and most importantly, we develop a computationally efficient estimation procedure using an information criterion derived from the minimal description length (MDL) principle, which substantially generalizes the procedure of [66] for TAR model. The procedure allows not only fast estimation of the number and the location of thresholds, but also model selection in each regime.

This paper is organized as follows. Section 2 defines the GTLVM. Section 3 establishes sufficient conditions for strict stationarity and ergodicity of GTLVM. Section 4 proposes a genetic algorithm using MDL principle for parameter estimation. Section 5 establishes asymptotic consistency and convergence rates of parameter estimations. Section 6 provides two simulation studies to illustrate the effectiveness of the proposed estimation method. Section 7 presents an application of modeling IPO volumes in U.S. stock market. The proof of Theorems are provided in the appendix.

## 2. Generalized threshold latent variable model

Consider a time series  $\{y_t\}_{t=1,\dots,n}$  in which  $y_t$  depends on its past observations and a latent variable  $\lambda_t$ , and follows different models when a threshold variable  $z_{t-d}$  belongs to different regimes. The latent variable  $\lambda_t$  is not only associated with past observations of  $y_t$  but also exogenous covariates  $X_t$ , and has a regime switching structure. Specifically, the conditional density function of  $y_t$  satisfies

$$f(y_t|Y_{t-1}, X_t, z_{t-d}) = \sum_{i=1}^{r+1} \left[ \int f_i(y_t|Y_{t-1}, \lambda_t) h_i(\lambda_t|Y_{t-1}, X_t) d\lambda_t \right] \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \tag{2.1}$$

where  $-\infty = \theta_0 < \theta_1 < \dots < \theta_r < \theta_{r+1} = \infty$  are the thresholds that classify  $y_t$  into  $r + 1$  regimes based on the threshold variable  $z_{t-d}$  and *threshold delay parameter*  $d$ ,  $Y_{t-1} = \{y_{t-1}, \dots, y_{t-p_Y}\}$  are the past observations, and  $X_t = \{x_{t,1}, \dots, x_{t,p_X}\}$  are the covariates. Moreover,  $f_i$  and  $h_i$  are conditional densities of  $y_t$  and  $\lambda_t$  of regime  $i$ , respectively. We assume that  $z_{t-d}$  is measurable with respect to the sigma-field generated by  $\{X_{t-d}, y_{t-d}, X_{t-d-1}, y_{t-d-1}, \dots\}$ . Typical examples include  $z_{t-d} = y_{t-d}$  for self-excited threshold models, and  $z_{k-d} = x_{k-d,k}$  for some exogenous covariates  $x_{t-d,k}$ .

To explicitly describe the dependence of  $y_t$  on  $\lambda_t$ ,  $Y_{t-1}$ ,  $X_t$  and  $z_{t-d}$ , the GTLVM can be specified as

$$\left\{ \begin{array}{l} y_t = \sum_{i=1}^{r+1} \left[ \omega_i + \sum_{l=1}^{p_{Y,1}} \psi_{y,1,i}^{(l)} g_{y,1}(y_{t-l}) + g_\lambda(\lambda_t, u_t) + \sum_{m=0}^{p_e} \psi_{\epsilon,i}^{(m)} \epsilon_{t-m} \right] \\ \quad \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \\ \phi(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y_{t-j}) + \sum_{k=1}^{p_X} \psi_{x,i}^{(k)} x_{t,k} + \sum_{q=0}^{p_e} \psi_{\epsilon,i}^{(q)} \epsilon_{t-q} \right] \\ \quad \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \end{array} \right. \tag{2.2}$$

where  $\phi(\cdot)$  is the link function which is smooth and monotonic increasing;  $\{\epsilon_t\}$  and  $\{e_t\}$  are independent and identical (i.i.d.) mean-zero errors;  $g_{y,1}(\cdot)$  and  $g_{y,2}(\cdot)$  are known continuous functions;  $g_\lambda(\lambda_t, u_t)$  is an inverse cumulative distribution function with parameter  $\lambda_t$ , and  $\{u_t\}$  are i.i.d. uniform random variables on  $[0, 1]$ . Without loss of generality, we may assume  $g_{y,1}(0) = g_{y,2}(0) = 0$  by adjusting  $\omega_i$  and  $\alpha_i$  accordingly. Note that (2.2) is a special case of (2.1): In (2.2), for any regime  $i$ , the conditional density of  $y_t$  given  $\lambda_t$ , denoted as  $f_i(y_t|Y_{t-1}, \lambda_t)$ , can be determined by the distributions of  $u_t$  and  $e_t$ . Also, the conditional density of  $\lambda_t$  given  $Y_{t-1}, X_t$ , denoted as  $h_i(\lambda_t|Y_{t-1}, X_t)$ , can be found from the second equation of (2.2). Integrating out the effect of  $\lambda_t$  in  $f_i(y_t|Y_{t-1}, \lambda_t)$  with respect to  $h_i$ , (2.1) follows. Although (2.1) is slightly more general than (2.2), we focus our attention on (2.2) for convenience in parametric modeling.

By properly choosing  $\lambda_t$ ,  $g_\lambda(\lambda_t, u_t)$ ,  $\phi(\lambda_t)$  and  $z_{t-d}$ , model (2.2) covers a number of the aforementioned threshold models in the literature. For instance, if  $g_\lambda(\lambda_t, u_t) = 0$  and  $g_{y,1}$  is the identity function, then (2.2) reduces to threshold autoregressive and moving average (TARMA) model

$$y_t = \sum_{i=1}^{r+1} \left[ \omega_i + \sum_{l=1}^{p_{Y,1}} \psi_{y,1,i}^{(l)} y_{t-l} + \sum_{m=0}^{p_e} \psi_{e,i}^{(m)} e_{t-m} \right] I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \quad (2.3)$$

as in [31]. If, in addition,  $p_e = 0$ , then (2.3) reduces to threshold autoregressive (TAR) model as in [58]. If  $g_\lambda(\lambda_t, u_t) = g_\lambda(u_t)$  is a quantile function with  $p_e = 0$ , then (2.2) reduces to the quantile self-excited threshold autoregressive (QSE-TAR) model in [8]. Next, denote  $\lambda_t = E(y_t)$  and let  $g_\lambda(\lambda_t, u_t) = F^{-1}(\lambda_t, u_t)$  be the inverse of the cumulative distribution function of some exponential family distribution with probability density

$$f(y_t; \lambda_t, a_t, \nu) = \exp \left[ \frac{1}{\nu a_t} \{y_t \eta(\lambda_t) - b(\lambda_t)\} + c(y_t; \nu a_t) \right],$$

where  $\eta(\lambda_t)$  is the canonical parameter,  $\nu$  is an overdispersion parameter and  $a_t$  is a user-specified weight. Then, the special case of (2.2) given by

$$\begin{cases} y_t = g_\lambda(\lambda_t, u_t), \\ \phi(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_Y} \psi_{y,i}^{(j)} g_y(y_{t-j}) + \sum_{k=1}^{p_X} \psi_{x,i}^{(k)} x_{t,k} + \psi_{\epsilon,i} \epsilon_t \right] \\ \quad \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \end{cases} \quad (2.4)$$

reduces to the GTMM in [53] when  $\phi(\lambda_t) = 0$  in the first or last regime, and reduces to the GTSRM in [54] when  $\psi_{\epsilon,j} = 0$  for all  $j$ . In addition, (2.2) covers models with double threshold structure and conditional heteroskedasticity. For example, with  $\lambda_t = \sigma_t$  and  $\phi(\sigma_t) = \sigma_t^2$ ,

$$\begin{cases} y_t = \sum_{i=1}^{r+1} \left[ \omega_i + \sum_{l=1}^p \psi_{y,1,i}^{(l)} y_{t-l} + \sigma_t \epsilon_t \right] I(y_{t-d} \in (\theta_{i-1}, \theta_i]), \\ \sigma_t^2 = \sum_{i=1}^{r+1} \left( \alpha_i + \sum_{j=1}^p \psi_{y,2,i}^{(j)} y_{t-j}^2 \right) I(y_{t-d} \in (\theta_{i-1}, \theta_i]), \end{cases} \quad (2.5)$$

is the MTDAR model ([33]). Similar arguments apply to the THSV model in [56].

Therefore, the proposed model (2.2) allows a unified treatment to the open problem of establishing the stationarity and ergodicity of many existing threshold-type models. Moreover, the algorithm developed in Section 4 provides an efficient solution to the computationally challenging problem of estimation and model selection for these models.

Denote the threshold parameters by  $\Theta = (\theta_1, \dots, \theta_r)$  and the model parameters for the  $i$ th regime by  $\Psi_i = (\omega_i, \Psi_{Y,1,i}, \Psi_{e,i}, \alpha_i, \Psi_{Y,2,i}, \Psi_{X,i}, \Psi_{\epsilon,i})$ , where  $\Psi_{Y,1,i} = (\psi_{y,1,i}^{(1)}, \dots, \psi_{y,1,i}^{(p_{Y,1,i})})$ ,  $\Psi_{e,i} = (\psi_{e,i}^{(0)}, \dots, \psi_{e,i}^{(p_e)})$ ,  $\Psi_{Y,2,i} = (\psi_{y,2,i}^{(1)}, \dots, \psi_{y,2,i}^{(p_{Y,2,i})})$ ,  $\Psi_{X,i} = (\psi_{x,i}^{(1)}, \dots, \psi_{x,i}^{(p_X)})$ ,  $\Psi_{\epsilon,i} = (\psi_{\epsilon,i}^{(0)}, \dots, \psi_{\epsilon,i}^{(p_\epsilon)})$ . Combining the parameters in all regimes, we define  $\Psi = (\Psi_1, \dots, \Psi_{r+1})$ . Note that we allow some of  $\psi_{y,1,i}^{(l)}$ ,  $\psi_{e,i}^{(m)}$ ,  $\psi_{y,2,i}^{(j)}$ ,  $\psi_{x,i}^{(k)}$  and  $\psi_{\epsilon,i}^{(a)}$  equal to zero so that different autoregressive orders and covariates can be included in different regimes. Denote the model order parameter as  $p = (p_1, \dots, p_{r+1})$ , where  $p_i = (p_{Y,1,i}, p_{e,i}, p_{Y,2,i}, p_{X,i}, p_{\epsilon,i})$ ,  $p_{Y,1,i}, p_{e,i}, p_{Y,2,i}$  and  $p_{\epsilon,i}$  are integers, and  $p_{X,i} = (p_{X,i}^{(1)}, \dots, p_{X,i}^{(p_X)})$  is a binary vector indicating the nonzero entries of  $\Psi_{X,i}$ . Thus, in (2.2),  $p_{Y,1} = \max_{i=1, \dots, r+1} p_{Y,1,i}$ , and  $p_{Y,2}, p_e$  and  $p_\epsilon$  are defined similarly.

### 3. Stationarity and ergodicity

We state a set of sufficient conditions for the strict stationarity and ergodicity of the GTLVM (2.2) as follows:

*Condition 1.*

a) The covariate  $X_t$  is independent of  $\{y_s\}_{s < t}$  and  $E|X_t| < \infty$ . In addition, there exists a positive integer  $\tilde{p}$  such that  $\{(X_t, \dots, X_{t-\tilde{p}+1})\}_{t=1, \dots}$  is Markovian, strictly stationary and ergodic. Moreover, there exists a non-negative integer  $q$  such that  $z_t$  is measurable with respect to the sigma-field generated by  $\{X_t, y_t, \dots, X_{t-q}, y_{t-q}\}$ .

b) The link function  $\phi(\cdot)$  is smooth and strictly increasing, and is either concave or a polynomial of order  $\gamma \geq 1$ ;

c) There exist constants  $b_1 > 0$  and  $\tilde{y} > 0$  such that  $|g_{y,1}(y)| \leq b_1|y|$  for all  $|y| > \tilde{y}$ , and  $g_{y,1}(y) \leq G_1$  for all  $|y| \leq \tilde{y}$ . In addition,  $\phi(\cdot)$  and  $g_{y,2}(\cdot)$  satisfy  $\phi^{-1}(g_{y,2}(y)) \leq b_2|y|$  for some  $b_2 > 0$  and all  $|y| > \tilde{y}$ , and  $\phi^{-1}(g_{y,2}(y)) \leq G_2$  for some constant  $G_2$  and all  $|y| \leq \tilde{y}$ ;

d)  $E[|g_\lambda(\lambda, u)| \mid \lambda]$  is increasing in  $\lambda$ . Moreover, there exist some positive constants  $b_\lambda, \tilde{\lambda}$  and  $H$  such that  $E[|g_\lambda(\lambda, u)| \mid \lambda] \leq b_\lambda|\lambda|$  if  $|\lambda| > \tilde{\lambda}$ , and  $E[|g_\lambda(\lambda, u)| \mid \lambda] \leq H$  if  $|\lambda| \leq \tilde{\lambda}$ ;

e)  $E(|e_t|) < \infty$ ,  $E[\phi^{-1}(\epsilon_t)] < \infty$ , and  $E[\phi^{-1}(x_{t,k})] < \infty$  for all  $x_{t,k}$ ;

f) For the case  $p_e > 0$ ,  $\{y_t\}$  is an irreducible process.

*Remark 1.* Denote  $\xi_t = (y_t, X_t, e_t, \epsilon_t)$  and  $p^* = \max\{p_{Y,1}, p_{Y,2}, \tilde{p}, p_e, p_\epsilon, q + d, 1\}$ . Condition 1a) ensures the Markovian property of  $\Xi_t = (\xi_t, \dots, \xi_{t-p^*+1})$ , which helps to prove the strict stationarity and ergodicity of  $\{y_t\}$ . Condition

1b) is satisfied for the concave link from the exponential family distribution (for example, log-link in [54]) and the square-link from modeling conditional heteroskedasticity ([33]). Conditions 1c) and 1d) are regularity conditions that control the increasing rates to avoid explosive behaviors. For example, taking  $\phi(\cdot)$  and  $g_{y,2}(\cdot)$  as the log-transform in [54], 1c) and 1d) hold with  $b_2 = b_\lambda = 1$ . Condition 1e) regulates the tails of the noise and covariates. For example, if  $\phi$  is the log link, then  $\epsilon_t$  must have tails lighter than Laplace distribution with intensity 1. Similar conditions are found in [11], [38] and [33]. Finally, the irreducibility required in Condition 1f) is used to derive the geometric ergodicity of the process when  $p_e > 0$ , see [37] and [38]. For  $p_e = 0$ , the irreducibility is guaranteed by the following lemma.

**Lemma 3.1.** *Under Condition 1, if  $p_e = 0$  in (2.2), then  $\{y_t\}$  is irreducible.*

To establish stationarity and ergodicity of a threshold model, a condition for preventing the series from exploding is required. For example,  $\max_i \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(j)}| < 1$  for TAR model in [3]. Under Condition 1, we have the following result for GTLVM, where the proof is provided in the Appendix.

**Theorem 3.1.** *Suppose that a process  $\{y_t\}$  satisfies (2.2) and Condition 1. Let*

$$\rho_i(\gamma) = \left( b_1 \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| + b_\lambda b_2 \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}|^{1/\gamma} \right).$$

*If either one of the following conditions holds:*

- 1)  $\phi$  is concave,  $\max_{i=1,\dots,r+1} \rho_i(1) < 1$ , and  $\max_{i=1,\dots,r+1} \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| < 1$ ;
- 2)  $\phi$  is a polynomial of order  $\gamma \geq 1$  and  $\max_{i=1,\dots,r+1} \rho_i(\gamma) < 1$ ,

*then  $\{y_t\}$  is strictly stationary and ergodic.*

*Remark 2.* The HAR model in [36] has a two-regime structure

$$\begin{aligned} y_t &= \begin{cases} \omega_1 + \sum_{l=1}^p \psi_{y,1}^{(l)} y_{t-l} + \sigma_1 e_t, & z_t^* = 1, \\ \omega_2 + \sum_{l=1}^p \psi_{y,2}^{(l)} y_{t-l} + \sigma_2 e_t, & z_t^* = 0, \end{cases} \\ z_t^* &= \begin{cases} 1, & \text{if } y_{t-d} \leq \theta_1, \\ 0, & \text{if } y_{t-d} > \theta_1 + a, \\ z_{t-1}^*, & \text{otherwise,} \end{cases} \end{aligned} \quad (3.1)$$

where  $a > 0$  and  $(\theta_1, \theta_1 + a]$  is called the hysteresis region, is covered by (2.2) by defining  $g_\lambda(\lambda_t, u_t) = 0$ ,  $p_e = 0$ , and

$$\begin{aligned} z_{t-d} &= y_{t-d} I(y_{t-d} \notin (\theta_1, \theta_1 + a]) \\ &+ \sum_{l=1}^{\infty} y_{t-d-l} I(y_{t-d-l} \notin (\theta_1, \theta_1 + a]) \prod_{j=0}^{l-1} I(y_{t-d-j} \in (\theta_1, \theta_1 + a]). \end{aligned} \quad (3.2)$$

Note that under (3.2), the threshold variable  $z_{t-d}$  is measurable with respect to the sigma field generated by  $\{y_{t-d}, y_{t-d-1}, \dots\}$  and thus does not satisfy Condition 1a). Nevertheless, Theorem 3.1 can be established if we replace Condition 1a) by:

Condition 1a') The covariate  $X_t$  is independent of  $\{y_s\}_{s < t}$  and  $E|X_t| < \infty$ . In addition, there exists a positive integer  $\tilde{p}$  such that  $\{(X_t, \dots, X_{t-\tilde{p}+1})\}_{t=1, \dots}$  is Markovian, strictly stationary and ergodic. Moreover, there exists some integer  $q$  such that the vector  $z_t^* = (I(z_{t-d} \in (\theta_0, \theta_1]), \dots, I(z_{t-d} \in (\theta_{r-1}, \theta_r]))$  is measurable with respect to the sigma-field generated by  $\{y_t, z_{t-1}^*, y_{t-1}, \dots, y_{t-q}, z_{t-q}^*\}$ . Also,  $z_t^*$  is irreducible and aperiodic.

Clearly, with  $q = 1$  and  $r = 1$ , Condition 1a') covers HAR model (3.1). In other words, Theorem 3.1 guarantees the stationarity and ergodicity of the multiple-regime extension of HAR model.

When specific knowledge is available on the threshold variable  $z_{k-d}$ , the conditions 1) and 2) in Theorem 3.1 can be relaxed as follows.

**Corollary 1.** *If  $z_{t-d} = x_{t-d,k}$ , then  $\{y_t\}$  is strictly stationary and ergodic under either one of the following conditions:*

- 1')  $\phi$  is concave,  $\sum_{i=1}^{r+1} \rho_i(1) \text{pr}(z_{t-d} \in (\theta_{i-1}, \theta_i]) < 1$  and  $\sum_{i=1}^{r+1} \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| \times \text{pr}(z_{t-d} \in (\theta_{i-1}, \theta_i]) < 1$ ;
- 2')  $\phi$  is polynomial of order  $\gamma \geq 1$  and  $\sum_{i=1}^{r+1} \rho_i(\gamma) \text{pr}(z_{t-d} \in (\theta_{i-1}, \theta_i]) < 1$ .

**Corollary 2.** *If  $z_{t-d} = y_{t-d}$ , then the conditions of  $\rho_i$  in Theorem 3.1 can be relaxed to  $\rho_1 < 1$ ,  $\rho_{r+1} < 1$  and  $\max\{\rho_1, \rho_{r+1}\}(1 - \pi_y) + \max_{i=2, \dots, r} \rho_i \pi_y < 1$ , where  $\pi_y = \sup_u \text{pr}(y_t \in (\theta_1, \theta_r] | Y_{t-1} = u)$ .*

We have the following corollary for the classical TAR and TARMA models defined in (2.3).

**Corollary 3.** *Suppose that  $E(|e_t|) < \infty$  and  $\max_{i=1, \dots, r+1} \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| < 1$ . If  $p_e = 0$  (TAR model), then (2.3) is stationary and ergodic. If  $p_e > 0$  (TARMA model), then (2.3) is stationary and ergodic provided that the irreducibility Condition 1f) holds.*

For TAR model, the condition in Corollary 3 is the same as [11]. For TARMA model, [38] also established the stationarity and ergodicity under the irreducibility condition. However, it requires  $\sum_{l=1}^{p_{Y,1}} \max_i |\psi_{y,1,i}^{(l)}| < 1$ , which is slightly stronger than the condition  $\max_i \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| < 1$  in Corollary 3. We remark that recently [14] proves the irreducibility for TARMA(1,1) model under some parametric conditions (see Condition (C3) in [14]). This justifies the potential validity of Condition 1f).

#### 4. Estimation and model selection criterion

Given the delay parameter  $d$ , the number of thresholds  $r$ , and the model order parameter  $p$ , the GTLVM is completely specified. Estimation of model parameters can be performed by maximum likelihood. Specifically, the log-likelihood

of the time series is

$$L(\Psi, \Theta, d) = \sum_{t=1}^n \left[ \sum_{i=1}^{r+1} l(\Psi_i; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right], \quad (4.1)$$

where  $l(\Psi_i; y_t, Y_{t-1}, X_t) = \log f(y_t | Y_{t-1}, X_t; z_{t-d} \in (\theta_{i-1}, \theta_i])$  is the conditional log-likelihood of  $y_t$  given  $\{Y_{t-1}, X_t\}$  in regime  $i$ . Let  $r^0, d^0, p^0, \Theta^0$  and  $\Psi^0$  be the true parameter values, and  $\hat{\Theta} = (\hat{\theta}_1, \dots, \hat{\theta}_r)$ ,  $\hat{\Psi} = (\hat{\Psi}_1, \dots, \hat{\Psi}_{r+1})$  be the corresponding maximum likelihood estimators that maximize (4.1) given  $d, r$  and  $p$ .

Based on the likelihood function, traditional methods such as sequential chi-square likelihood ratio test from [12] can be derived to determine the number of thresholds. However, different autoregressive orders and combination of covariates in different regimes contribute to complication in implementing the traditional methods. To overcome the computational burden, we adopt a model selection approach by developing a criterion function based on the minimal description length (MDL) principle ([47, 48]). Given a model  $\mathcal{M}$ , the criterion is defined by

$$\text{MDL}(\mathcal{M}) = \text{CL}(\mathcal{M}) + \text{CL}(\mathcal{E}_n | \mathcal{M}), \quad (4.2)$$

where the right hand side of (4.2) are the code lengths in bits for encoding the model and the fitted residuals  $\mathcal{E}_n = (\hat{e}_1, \dots, \hat{e}_n)$  given the model, respectively. Encoding the model  $\mathcal{M}$  requires the specification of  $r, d, p, \hat{\theta}_i$ s and  $\hat{\Psi}_i$ s. Thus,  $\text{CL}(\mathcal{M})$  can be expressed as

$$\text{CL}(\mathcal{M}) = \text{CL}(r) + \text{CL}(d) + \text{CL}(\hat{\theta}_1, \dots, \hat{\theta}_n) + \text{CL}(\hat{\Psi}_1) + \dots + \text{CL}(\hat{\Psi}_{r+1}).$$

From [47, 48] and [30], it requires approximately  $\log_2(n)$  bits to encode an integer and  $(\log_2 n)/2$  bits to encode a maximum likelihood estimator with  $n$  data points. From [66], the thresholds can be associated with the order statistics of the threshold variables  $\{z_1, \dots, z_n\}$  and require  $\sum_{i=1}^r \log_2(n_i)/2$  bits, where  $n_i$  is the number of observations in the  $i$ th regime. Recall that  $\psi_{x,i}^{(k)} = 0$  when the  $k$ th covariate is not included in the model. As the integer 0 requires 1 bit and  $\log_2 1 = 0$ , by denoting  $p'_{X,i} = \sum_{w=1}^{p_X} p_{X,i}^{(w)}$  as the number of nonzero entries in  $p_{X,i}$ , the maximum likelihood estimator  $\hat{\Psi}_{X,i}$  can be encoded with  $(p'_{X,i}/2) \log_2(n_i)$  bits. Similar arguments suggest that encoding  $\hat{\Psi}_{Y,1,i}, \hat{\Psi}_{e,i}, \hat{\Psi}_{Y,2,i}$  and  $\hat{\Psi}_{\epsilon,i}$  require  $(p_{Y,1,i}/2) \log_2(n_i), (p_{e,i}/2) \log_2(n_i), (p_{Y,2,i}/2) \log_2(n_i)$  and  $(p_{\epsilon,i}/2) \log_2(n_i)$  bits, respectively. Thus, we have

$$\begin{aligned} \text{CL}(\mathcal{M}) &= \log_2(r) + \log_2(d) + \sum_{i=1}^r \frac{\log_2(n_i)}{2} + \sum_{i=1}^{r+1} \log_2(p'_i + 4) \\ &\quad + \sum_{i=1}^{r+1} \frac{p'_i + 4}{2} \log_2(n_i), \end{aligned}$$

where  $p'_i = p_{Y,1,i} + p_{e,i} + p_{Y,2,i} + p'_{X,i} + p_{\epsilon,i}$  is the total number of nonzero coefficients in regime  $i$ . From [47],  $\text{CL}(\mathcal{E}_i | \mathcal{M})$  can be approximated by the negative



of  $\log_2$  of the likelihood. Hence,

$$\begin{aligned} \text{MDL}(\mathcal{M}) &= \log_2(r) + \log_2(d) + \sum_{i=1}^r \frac{\log_2(n_i)}{2} + \sum_{i=1}^{r+1} \log_2(p'_i + 4) \\ &\quad + \sum_{i=1}^{r+1} \frac{p'_i + 4}{2} \log_2(n_i) \\ &\quad - \log_2(e)L(\hat{\Psi}, \hat{\Theta}, d). \end{aligned} \tag{4.3}$$

The optimal model can then be selected as the values  $\hat{r}$ ,  $\hat{d}$  and  $\hat{p}$  that minimize the  $\text{MDL}(\mathcal{M})$ .

Since the likelihood function remains constant when  $\theta_i \in (R_j, R_{j+1}]$ ,  $j = 1, \dots, n-1$ , where  $\{R_1, \dots, R_n\}$  is the ordered observations of the threshold variable, the estimator  $\hat{\theta}_i$  may take any value on  $(R_j, R_{j+1}]$  for some  $j$ . Without loss of generality we take  $\hat{\theta}_i = R_{j+1}$ . Then, to obtain an approximate solution to the optimization problem, we developed a genetic algorithm which is found to achieve promising performance for related optimization problems in change-point analysis ([18, 20], [41]) and estimation of the TAR model ([64] and [66]). Inspired by [66], we develop the methodology with modifications for a simultaneous detection of both autoregressive and covariate structures.

The genetic algorithm is an imitation of the biological evolution for optimization. It involves inheritance, crossover, mutation and filtering. Specifically, the algorithm begins with a population of *chromosomes*, where each chromosome stores the information of a candidate solution to the optimization problem. For our application, we first fix a  $d$  and define each candidate solution as a model  $\mathcal{M}$  specified by some  $r$  and  $p$ . Then, the performance of each chromosome is measured by its information criterion  $\text{MDL}(\mathcal{M})$ . Chromosomes with better performance have higher probabilities to conduct crossover and produce *offspring*, and so their good model features are more likely to be inherited. Meanwhile, mutation occurs with a small probability, which brings in new models to seek the global optimum. After several generations of crossover and mutation, the best performing model is selected as the optimal one. Finally, we repeat the procedure with different  $d$  and select the best performing model that attains the smallest  $\text{MDL}(\mathcal{M})$ .

Specifically, each step of the genetic algorithm is described as follows.

**Chromosome Formation:** First, generate the initial *population*, which is a set of chromosomes in vector form. Each chromosome is expressed as  $c = [r, p_1, (\theta_1, p_2), \dots, (\theta_r, p_{r+1})]$ ; the parameter estimate  $\hat{\Psi}$  is obtained once  $c$  is specified. Similar to [66], the initial population is created as follows:

- 1) The number of thresholds  $r$  is generated from a Poisson distribution with mean 2.
- 2) Sample  $\theta_i$ s uniformly from  $\{z_i\}$ . Reject the sample and sample again if any regime has fewer than  $\tau_{n,0}$  observations. This *minimum span condition* ensures the estimation accuracy of  $\Psi$ .

- 3) Select  $p_{Y,1,i}$ ,  $p_{e,i}$ ,  $p_{Y,2,i}$  and  $p_{\epsilon,i}$  uniformly from  $\{0, \dots, Q_{Y,1}\}$ ,  $\{0, \dots, Q_e\}$ ,  $\{0, \dots, Q_{Y,2}\}$  and  $\{0, \dots, Q_\epsilon\}$ , respectively, where  $Q_{Y,1}$ ,  $Q_e$ ,  $Q_{Y,2}$  and  $Q_\epsilon$  are pre-specified upper bounds of model orders. Generate a binary vector  $p_{X,i}$  of length  $p_X$  with each element following an independent Bernoulli distribution with mean 0.5. Set  $p_i = (p_{Y,1,i}, p_{e,i}, p_{Y,2,i}, p_{X,i}, p_{\epsilon,i})$ .

The  $\text{MDL}(\mathcal{M})$  of each chromosome is then computed by (4.3).

**Crossover and Mutation:** Crossover and mutation are two methods for generating offspring. In crossover, two chromosomes are selected from the population as “parents” with probabilities proportional to the inverse of their ranks of  $\text{MDL}(\mathcal{M})$ . Next, a  $p_1^o$  is drawn from one of the parent’s  $p_1$  with equal probability. Then, for both parents, each of their  $\{\theta_j, p_{j+1}\}$  is selected with probability 0.5. Sort all selected  $\{\theta_j, p_{j+1}\}$ s in ascending order of  $\theta_j$  to produce an offspring  $[r^o, p_1^o, (\theta_1^o, p_2^o), \dots, (\theta_r^o, p_{r^o+1}^o)]$ . If some thresholds  $\theta_j^o$ s violate the minimum span condition, randomly delete the pair  $\{\theta_j^o, p_{j+1}^o\}$  until the condition is satisfied.

In mutation, one parent chromosome is selected from the population with probabilities proportional to the inverse of the ranks of  $\text{MDL}(\mathcal{M})$ . Then, a new chromosome is generated to crossover with the parent chromosome to produce an offspring. To achieve a higher degree of exploration, the features from the generated parent are selected with probability 0.7.

To balance between retaining good features with crossover and bringing new solutions with mutation, the probabilities of performing crossover and mutation are 0.9 and 0.1, respectively. To explore more possibilities in the model order, every offspring will have its order parameter in one randomly selected regime replaced by a newly generated order, with probability 0.3.

After conducting crossover and mutation, the group of offspring become a new generation of chromosomes. To ensure monotonicity of optimization, an *elitist* step is conducted to replace the worst performing 20 chromosomes in the new generation by the best performing chromosome in the previous generation.

**Migration:** With the advance of parallel computing, the *island model* is introduced to accelerate the computation and alleviate trapping in suboptimal solutions. In particular, we perform genetic algorithm on  $N_I$  groups of subpopulations with size  $N_p$ . These subpopulations conduct their own reproduction steps and thus are treated as distinct islands. To share good features between islands, for every  $M_i$  generations, the  $M_N$  worst performing chromosomes in the  $j$ th island are replaced by the best  $M_N$  chromosomes from the  $(j - 1)$ th island, for  $j = 1, \dots, N_I$ , where the 0th island is conventionally defined as the  $N_I$ -th island. In this article we used  $N_I = 50$ ,  $N_p = 200$ ,  $M_i = 4$ , and  $M_N = 2$ . The full mechanism and improvement in performance of the parallelized genetic algorithm can be found in [1, 2].

**Claim of Convergence:** When the best chromosome remains unchanged over 20 generations, we claim that convergence is achieved and the optimal model is obtained from the parameters of the best chromosome. Alternatively, in consideration of computational efficiency, the algorithm may be stopped after a fixed number of generations.

## 5. Asymptotic inferences of GTLVM

### 5.1. Assumptions for asymptotic inferences

Apart from Condition 1, we state the following assumptions for asymptotic inferences. First, Assumptions 1–3 are proposed for the consistency:

*Assumption 1.* The parameters  $(\Psi, \Theta, d)$  are in the space  $\Omega_{\mathcal{M}} \times \{1, \dots, D\}$ , where  $\Omega_{\mathcal{M}} = \Omega_{\Psi} \times \Omega_{\Theta}$  is a compact subset of  $\mathbb{R}^{(p_{Y,1}+p_e+p_{Y,2}+p_X+p_e+4) \times (r+1)} \times \mathbb{R}^r$ , and  $D$  is some positive integer. In addition,  $(\Psi^0, \Theta^0)$  is an interior point in  $\Omega_{\mathcal{M}}$ .

*Assumption 2.* The conditional density  $f(y_t | Y_{t-1}, X_t, z_{t-d})$  is regular in the sense that the maximum likelihood estimator  $\hat{\Psi}_i$  is asymptotically normal ([17] and [44]). We assume  $E[\{\partial \log(f(y_t | Y_{t-1}, X_t, z_{t-d}))/\partial \Psi_i\}^2] < \infty$ . Moreover,  $l(\Psi_i; y_t, Y_{t-1}, X_t)$  is concave in  $\Psi_i$ , and  $\partial^3 l(\Psi_i; y_t, Y_{t-1}, X_t)/\partial \Psi_i^3$  exists for  $i = 1, \dots, r + 1$  and is bounded by an integrable function in the neighborhood of  $\Psi_i^0$ . Furthermore,  $f(y_t | Y_{t-1}, X_t, z_{t-d} \in (\theta_{i-1}^0, \theta_i^0))$  and  $f(y_t | Y_{t-1}, X_t, z_{t-d} \in (\theta_i^0, \theta_{i+1}^0])$  are not equal almost everywhere for all  $i$  and all  $\{Y_{t-1}, X_t\}$ .

*Assumption 3.* The joint probability density of  $\{z_{t-i}, z_{t-j}\}$ ,  $\pi_{z,i-j}(\cdot, \cdot)$ , is uniformly bounded. In addition, for any vector  $\Phi$  that has the same dimension as  $(1, X_t^T)$  and satisfies  $|\Phi| = 1$ , there exists an  $\epsilon > 0$  such that

$$\text{pr}(|\Phi(1, X_t^T)^T| > \epsilon | z_{t-i}, z_{t-j}) > 0 \quad \text{almost surely,} \quad (5.1)$$

with respect to the joint distribution of  $(z_{t-i}, z_{t-j})$ , where  $i, j = 1, \dots, D$ .

Assumptions 1–2 are standard regulatory conditions for statistical inference in parametric models. The assumption on the third-order derivatives of  $L(\Psi, \Theta, d)$  is essential in deriving the asymptotic distribution of  $\hat{\Theta}$ ; see Theorem 5.41 of [62]. Assumption 3, which is in similar spirit as Assumption 3 in [54], assumes the linear independence of  $X_t$  to eliminate redundancy in the covariates. It holds if the joint conditional density of the exogenous covariates  $X_t$  is non-degenerate given  $z_{t-i}$  and  $z_{t-j}$ .

Additionally, the following two assumptions are required for the convergence rates of estimators:

*Assumption 4.* There exists an integrable function  $\Gamma(Y_{t-1}, X_t, y_t)$  satisfying

$$\left| l(\Psi^{(1)}; y_t, Y_{t-1}, X_t) - l(\Psi^{(2)}; y_t, Y_{t-1}, X_t) \right| < \Gamma(Y_{t-1}, X_t, y_t) |\Psi^{(1)} - \Psi^{(2)}| \quad \text{a.s.,}$$

for any  $\Psi^{(1)}, \Psi^{(2)} \in \Omega_{\Psi}$ . Furthermore, we assume that either one of the following condition holds:

- a) The marginal density of  $z_t$ ,  $\pi_z(\cdot)$ , is continuous at  $\{\theta_1^0, \dots, \theta_r^0\}$ , and the joint probability density of  $\{z_{t-i}, z_{t-j}\}$ ,  $\pi_{z,i-j}(\cdot, \cdot)$ , is positive everywhere. Moreover, there exists an  $\epsilon > 0$  such that, for all  $z_{t-d^0} \in [\theta_i^0 - \epsilon, \theta_i^0 + \epsilon]$ ,  $i = 1, \dots, r$ ,  $E[\Gamma^2(Y_{t-1}, X_t, y_t)] < \infty$ . The joint conditional distribution function of  $\{Y_{t-1}, X_t\}$  given  $z_{t-d^0}$  is continuous at  $z_{t-d^0} = \theta_1^0, \dots, \theta_r^0$ .

- b) The threshold variable  $z_t$  is discrete, and  $y_t$  is in the exponential family satisfying

$$\begin{cases} f(y_t; \lambda_t, a_t) = \sum_{i=1}^{r+1} \exp \left[ \frac{1}{\nu_i a_t} \{T(y_t)\gamma(\lambda_t) - b(\lambda_t)\} + c(y_t; \nu_i a_t) \right] \\ \quad \times I(z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]), \\ \phi(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_Y} \psi_{y,i}^{(j)} g_y(y_{t-j}) + \sum_{k=1}^{p_X} \psi_{x,i}^{(k)} x_{t,k} \right] \\ \quad \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \end{cases} \quad (5.2)$$

where  $T(\cdot)$  is a measurable function such that  $\text{var}(T(y_t) \mid \lambda_t) \in (0, \infty)$  almost surely,  $\gamma(\cdot), b(\cdot)$  are continuous transformation functions, and  $\nu_i$  is the overdispersion parameter for regime  $i$ . Moreover, we assume that  $\mathbb{E}[\Gamma^2(Y_{t-1}, X_t, y_t)] < \infty$ ,  $\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n c_t^{3/2} < \infty$  and  $\log[\mathbb{E}(e^{uT(y_t)} \mid \lambda_t)] \leq c_t$  for some constant  $u$  and sequence  $\{c_t\}$ .

*Assumption 5.* When  $z_t$  is continuous, there exists some  $\Delta > 0$  such that the process  $\{\Gamma(Y_{t-1}, X_t, y_t)I(z_{t-d^0} \in [\theta_i^0 - \Delta, \theta_i^0 + \Delta]), z_{t-d^0}I(z_{t-d^0} \in [\theta_i^0 - \Delta, \theta_i^0 + \Delta])\}$  is  $\rho$ -mixing with summable mixing coefficients.

Assumption 4 are the conditions for the likelihood function with respect to continuous and discrete threshold variables. Asymptotic theory for threshold models with discrete threshold variables does not seem to have been studied in the literature. Similar to Assumption 6 in [54], we impose a square-integrable bound function for the difference of log-likelihood in both Assumption 4a) and 4b). While an example of verifying Assumption 4a) has been shown in supplementary materials of [54], an example of verifying Assumption 4b) is given in the Appendix. Furthermore, Assumption 4b) requires that the conditional density of  $y_t$  given  $\{Y_{t-1}, X_t, z_{t-d}\}$  is in the exponential family. Thus, an explicit form of the difference in log-likelihood is available for applying results in large deviation theory; see [46]. Note that although the assumptions and asymptotic properties of  $\hat{\Theta}$  for continuous and discrete  $z_t$  are different, the same estimation procedure proposed in Section 4 is applied. Assumption 5 is analogous to Assumption 7 in [54] for proving the convergence of  $n(\hat{\Theta} - \Theta^0)$ . For any integer  $j$ , let  $\mathcal{A}$  and  $\mathcal{A}^*$  be the  $\sigma$ -algebras generated by  $\{w_t\}_{t \leq j}$  and  $\{w_t\}_{t \geq j+k}$ , respectively. If  $\{w_t\}$  is  $\rho$ -mixing, then there exists a sequence  $\{\rho(k)\}_{k=1,2,\dots}$  with  $\lim_{k \rightarrow \infty} \rho(k) \rightarrow 0$  such that, for all square-integrable random variables  $g$  and  $h$  that are respectively  $\mathcal{A}$  and  $\mathcal{A}^*$ -measurable,  $|\text{corr}(g, h)| \leq \rho(k)$  holds; see [7] and [21]. See also Examples 1 and 2 in the supplementary materials of [54] for verification of  $\rho$ -mixing and selections of the function  $\Gamma$ . We will illustrate the verification of the above assumptions in some explicit examples in Section 6.

## 5.2. Asymptotic theorems

Under Condition 1, we develop the consistency and asymptotic properties of the parameter estimates. The proofs are provided in the appendix.

**Theorem 5.1.** *If Assumptions 1–3 hold, then  $\hat{d} \rightarrow d^0$  and  $\hat{r} \rightarrow r^0$  almost surely. In addition, on  $\{\hat{d} = d^0, \hat{r} = r^0\}$ ,  $\hat{\Theta} \rightarrow \Theta^0$ ,  $\hat{p} \rightarrow p^0$  and  $\hat{\Psi} \rightarrow \Psi^0$  almost surely.*

Next, we derive the convergence rate of  $\hat{\Theta}$  for continuous or discrete  $z_t$ , respectively:

**Theorem 5.2.** *Under Assumptions 1–3,*

- 1) *if Assumptions 4a) and 5 are satisfied, then  $|\hat{\Theta} - \Theta^0| = O_p(n^{-1})$ ;*
- 2) *if Assumption 4b) is satisfied, then  $\text{pr}(\hat{\Theta} = \Theta^0) = 1 - O(n^{1/2}e^{-an})$  for some  $a > 0$ .*

Denote the difference between the log-likelihood of  $y_s$  under parameters  $\Psi^{(1)}$  and  $\Psi^{(2)}$  by

$$\zeta_s(y_s; Y_{s-1}, X_s, \Psi^{(1)}, \Psi^{(2)}) = l(\Psi^{(1)}; y_s, Y_{s-1}, X_s) - l(\Psi^{(2)}; y_s, Y_{s-1}, X_s).$$

Define a double-sided compound Poisson process  $\tilde{\ell}_i^*(\kappa_i) = \tilde{\ell}_{1,i}^*(\kappa_i)I(\kappa_i \geq 0) + \tilde{\ell}_{2,i}^*(\kappa_i)I(\kappa_i < 0)$ , where

$$\tilde{\ell}_{1,i}^*(\kappa_i) = \sum_{s=1}^{N_{1,i}(\kappa_i)} \zeta_s(y_s^*; Y_{s-1}^*, X_s^*, \Psi_{i+1}^0, \Psi_i^0), \tag{5.3}$$

is a compound Poisson process such that  $N_{1,i}(\kappa_i)$  is a Poisson processes with intensity  $\pi_z(\theta_i^0)$ , and  $(y_s^*, Y_{s-1}^*, X_s^*)$  is an independent copy of  $(y_s, Y_{s-1}, X_s)$  given  $z_{s-d} = (\theta_i^0)^+$ . Moreover,

$$\tilde{\ell}_{2,i}^*(\kappa_i) = \sum_{s=1}^{N_{2,i}(-\kappa_i)} \zeta_s(y_s^*; Y_{s-1}^*, X_s^*, \Psi_i^0, \Psi_{i+1}^0), \tag{5.4}$$

is a compound Poisson process, independent of  $\tilde{\ell}_{1,i}^*(\kappa_i)$ , defined by where  $N_{2,i}(-\kappa_i)$  is a Poisson processes with the same intensity  $\pi_z(\theta_i^0)$ , and  $(y_s^*, Y_{s-1}^*, X_s^*)$  is an independent copy of  $(y_s, Y_{s-1}, X_s)$  given  $z_{s-d} = (\theta_i^0)^-$ . In addition, define an aggregated compound Poisson process  $\tilde{\ell}^*(\kappa) = \sum_{i=1}^{r+1} \tilde{\ell}_i^*(\kappa_i)$ . The following theorems derive the asymptotic distribution of the threshold parameters  $\hat{\Theta}$  and model parameters  $\hat{\Psi}$ .

**Theorem 5.3.** *Under Assumptions 1–3, 4a) and 5, then  $n(\hat{\Theta} - \Theta^0)$  weakly converges to  $M^-$ , where the random  $r$ -dimension cube  $[M^-, M^+] = [M_1^-, M_1^+] \times [M_2^-, M_2^+] \times \dots \times [M_r^-, M_r^+]$  is an almost surely minimizer of the compound Poisson process  $\tilde{\ell}^*(\kappa)$ . In addition, if  $f(y | Y_{t-1}, X_t, z_{t-d})$  is continuous in  $y$ , then  $[M^-, M^+]$  is unique.*

*Remark 3.* Note that  $[M^-, M^+]$  depends on  $\Theta^0$  and  $\Psi^0$ . For constructing confidence intervals for  $\hat{\Theta}$ , algorithms for estimating  $M^-$  can be derived by similar methods in [32] and [67].

**Theorem 5.4.** *Define*

$$L'(\Psi, \Theta, d) = \frac{\partial}{\partial \Psi} \left[ \frac{1}{n} L(\Psi, \Theta, d) \right], \quad L''(\Psi, \Theta, d) = \frac{\partial^2}{\partial \Psi^2} \left[ \frac{1}{n} L(\Psi, \Theta, d) \right].$$

Under Assumptions 1–5,  $n^{1/2}(\hat{\Psi} - \Psi^0) \rightarrow_d N(0, \Sigma^*)$ , where

$$\Sigma^* = - [\mathbb{E}\{L''(\Psi^0, \Theta^0, d^0)\}]^{-1}. \quad (5.5)$$

Moreover,  $\Sigma^*$  is a diagonal block matrix, and thus  $\{\hat{\Psi}_i\}_{i=1, \dots, r+1}$  are asymptotically independent. In addition,  $n(\hat{\Theta} - \Theta^0)$  and  $n^{1/2}(\hat{\Psi} - \Psi^0)$  are asymptotically independent.

## 6. Simulation

In this section, two simulation experiments are performed to illustrate model fitting with respect to discrete and continuous thresholds. In each simulation experiment, 300 replications are conducted for every scenario. For simplicity,  $d$  is assumed to be known.

### 6.1. Example 1

Consider a three-regime self-excited GTLVM with Poisson distribution and log-link:

$$\begin{cases} f(y_t; \lambda_t) &= \lambda_t^{y_t} e^{-\lambda_t} / y_t!, \\ \log(\lambda_t) &= \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,i}^{(j)} \log(y_{t-j} + 1) \right] I(y_{t-4} \in (\theta_{i-1}, \theta_i]), \end{cases} \quad (6.1)$$

with  $\{\theta_1, \theta_2\} = \{24, 42\}$ ,  $\{\alpha_1, \alpha_2, \alpha_3\} = \{0.45, 3.4, 6.3\}$ ,  $\{p_{Y,2,1}, p_{Y,2,2}, p_{Y,2,3}\} = \{2, 2, 1\}$ , and  $\{(\psi_{y,1}^{(1)}, \psi_{y,1}^{(2)}), (\psi_{y,2}^{(1)}, \psi_{y,2}^{(2)}), \psi_{y,3}^{(1)}\} = \{(0.65, 0.25), (0.4, -0.35), -0.95\}$ . A time series plot of realization is shown in Figure 1.

Here we verify Condition 1, and the assumptions in Section 5 for model (6.1). Note that (6.1) can be expressed as (2.2) with  $\phi(x) = \log(x)$ ,  $g_\lambda(\lambda_t, \cdot)$  being the inverse c.d.f. of Poisson distribution with parameter  $\lambda_t$ ,  $g_{y,1}(x) \equiv 0$ ,  $g_{y,2}(x) = \log(x + 1)$ ,  $p_X = p_e = p_\epsilon = 0$ , and  $\psi_{y,1,i}^l = 0$  for all  $i$  and  $l$ . As no covariate is used and  $p_e = p_\epsilon = 0$ , it suffices to check Condition 1 b), c), d). From the forms of  $\phi$ ,  $g_{y,1}$  and  $g_{y,2}$ , Conditions 1b) and c) clearly hold with  $b_1 = b_2 = 1.1$ . Given  $\lambda$  and the uniform random variable  $u$  on  $[0, 1]$ ,  $g_\lambda(\lambda, u)$  is a Poisson random variable with parameter  $\lambda$ . Thus,  $\mathbb{E}[g_\lambda(\lambda, u) \mid \lambda] = \lambda$  and Condition 1d) holds with  $b_\lambda = 1$ . Therefore, Condition 1 holds and Theorem 3.1 implies that (6.1) is stationary and ergodic.

As the threshold variable is discrete and no covariate is used, it suffices to verify Assumptions 1, 2 and 4b). As is common in the literature, Assumption 1 can be achieved by focusing attention on a sufficiently large compact subset of the parameter space. For Assumption 2, note that the conditional density  $f(y_t \mid Y_{t-1}, y_{t-4})$  is the density of Poisson distribution, and is thus regular. Also,  $\mathbb{E}\{[\partial \log(f(y_t \mid Y_{t-1}, y_{t-4})) / \partial \psi_{y,i}^{(j)}]^2\} = E((y_t - \lambda_t)(\log(y_{t-j} + 1))^2) < \infty$ . Moreover, the third-order derivative of  $l(\Psi, \Theta, d)$  with respect to  $\psi_{y,i}^{(j)}$  exists and is equal to  $-\lambda_t(\log(y_{t-j} + 1))^3$ . Also,  $f(y_t \mid Y_{t-1}, y_{t-4} \in (\theta_{i-1}^0, \theta_i^0])$  and

$f(y_t | Y_{t-1}, y_{t-4} \in (\theta_i^0, \theta_{i+1}^0])$  are not equal almost everywhere for all  $i$  since the latent variable  $\lambda_t$  takes different values in the two regimes. Verification of Assumption 4b) is more technical and is given in the Appendix.

TABLE 1

Example 1: Percentage of correct number of estimated thresholds and correct model structure specification. MDL: minimum description length; BIC: Bayesian information criterion.

Size	Regime structure classified(%), MDL	Model structure classified(%), MDL	Regime structure classified(%), BIC	Model structure classified(%), BIC
200	91.7	65.3	85.6	60.0
500	96.7	94.0	94.3	87.3
1000	99.3	96.3	98.3	94.3
2000	100.0	97.0	99.3	93.7

TABLE 2

Example 1: Mean and standard derivations of parameter estimates. Parentheses: empirical standard derivations. The threshold estimates are based on replications with correct number of regimes. The coefficient estimates are based on replications with correct model structure.

Size	$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\alpha}_1$	$\hat{\psi}_{Y,2,1}^{(1)}$	$\hat{\psi}_{Y,2,1}^{(2)}$	$\hat{\alpha}_2$	$\hat{\psi}_{Y,2,2}^{(1)}$	$\hat{\psi}_{Y,2,2}^{(2)}$	$\hat{\alpha}_3$	$\hat{\psi}_{Y,2,3}^{(1)}$
200	24.036 (0.875)	42.007 (0.171)	0.430 (0.324)	0.631 (0.091)	0.276 (0.068)	3.411 (0.234)	0.395 (0.050)	-0.348 (0.052)	6.260 (0.267)	-0.938 (0.083)
500	24.003 (0.228)	42.000 (0.000)	0.463 (0.214)	0.647 (0.050)	0.250 (0.043)	3.396 (0.126)	0.397 (0.033)	-0.346 (0.032)	6.285 (0.131)	-0.946 (0.041)
1000	24.006 (0.082)	42.000 (0.000)	0.452 (0.174)	0.649 (0.039)	0.250 (0.034)	3.400 (0.094)	0.399 (0.023)	-0.349 (0.023)	6.293 (0.097)	-0.948 (0.030)
2000	24.000 (0.000)	42.000 (0.000)	0.447 (0.102)	0.652 (0.025)	0.249 (0.021)	3.402 (0.070)	0.401 (0.016)	-0.351 (0.015)	6.293 (0.063)	-0.948 (0.020)
True	24	42	0.45	0.65	0.25	3.4	0.4	-0.35	6.3	-0.95

Table 1 reports the model selection performance. Even for a small sample of size of 200, the percentage of correct number of estimated threshold is over 80%, and the percentage of correct identification of model order in all regimes is over 60%. For comparison, we repeated the experiment with the Bayesian Information Criterion, which is defined as  $-2L(\Psi, \Theta, d) + \sum_{i=1}^{r+1} \log(n_i)p'_i$ . It is found that minimal description length gives better performance. Other information criteria such as NAIC in [59] might also be adopted; however, NAIC and other AIC-type criteria are not consistent in estimating the true order of the model. See [26] for details about consistency of information criterion.

Furthermore, Table 2 summarizes the performance of thresholds and model parameters estimates within the replications that correctly specify the regime and model structure using minimal description length. A fast convergence speed for discrete thresholds is observed. Furthermore,  $O_p(n^{-1/2})$  convergence rate of  $\hat{\Psi}$  is realized based on the rapid convergence of  $\hat{\Theta}$ .

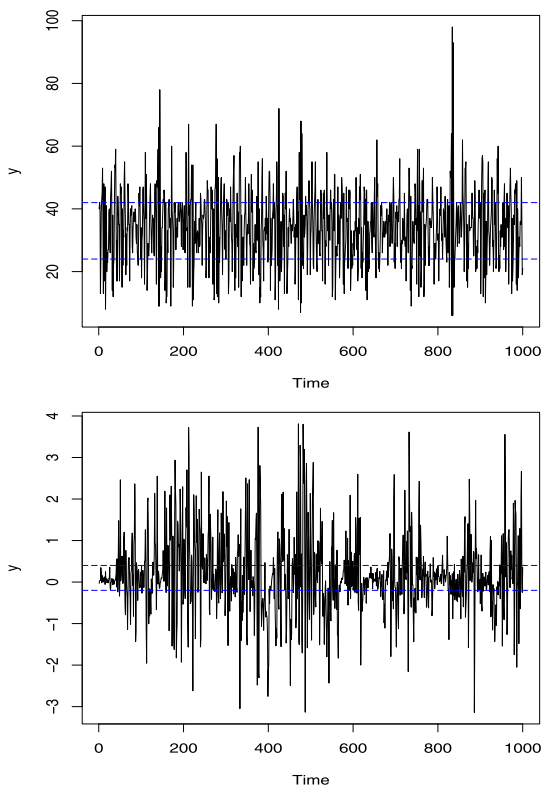


FIG 1. Sample plots of simulated series. Left: Example 1, with thresholds at 24 and 42 (horizontal dashed line). Right: Example 2, with thresholds at  $-0.2$  and  $0.4$  (horizontal dashed line)

## 6.2. Example 2

We consider a double threshold autoregressive model with conditional heteroskedasticity by a log-link on  $\sigma_t^2$ :

$$\begin{cases} y_t &= \sum_{i=1}^{r+1} \left[ \omega_i + \sum_{l=1}^{p_{Y,1}} \psi_{y,1,i}^{(l)} y_{t-l} + \sigma_t \epsilon_t \right] I(y_{t-4} \in (\theta_{i-1}, \theta_i]), \\ \log(\sigma_t^2) &= \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} \log(y_{t-j}^2) \right] I(y_{t-4} \in (\theta_{i-1}, \theta_i]). \end{cases} \quad (6.2)$$

Here  $\{\theta_1, \theta_2\} = \{-0.2, 0.4\}$ ,  $\{\omega_1, \omega_2, \omega_3\} = \{-0.15, 0.1, 0.3\}$ ,  $\{p_{Y,1,1}, p_{Y,1,2}, p_{Y,1,3}\} = \{1, 1, 2\}$ ,  $\{(\psi_{y,1,1}^{(1)}), (\psi_{y,1,2}^{(1)}), (\psi_{y,1,3}^{(1)}), \psi_{y,1,3}^{(2)}\} = \{(0.6), (0.25), (0.25, -0.7)\}$ ,  $\{\alpha_1, \alpha_2, \alpha_3\} = \{-0.4, -0.2, 0.15\}$ ,  $\{p_{Y,2,1}, p_{Y,2,2}, p_{Y,2,3}\} = \{1, 2, 0\}$ ,  $\{(\psi_{y,2,1}^{(1)}), (\psi_{y,2,2}^{(1)}), \psi_{y,2,2}^{(2)}), (\psi_{y,2,3}^{(1)})\} = \{(0.35), (0.45, 0.25), (0)\}$ , with  $\epsilon_t \stackrel{iid.}{\sim} \mathcal{N}(0, 1)$ . In the third regime, we select  $p_{Y,2,3} = 0$  so  $\sigma_t^2 = e^{0.15}$  is a constant. A sample time series plot is depicted in Figure 1.



Note that (6.2) can be expressed as (2.2) with  $\phi(x) = \log(x)$ ,  $g_\lambda(\lambda_t, \cdot) \equiv 0$ ,  $g_{y,1}(x) = x$ ,  $g_{y,2}(x) = 2 \log(x)$  and  $p_X = p_e = p_\epsilon = 0$ . Condition 1 can be verified readily using similar arguments in the verification for model (6.1). As the threshold variable is continuous and no covariate is used, it suffices to verify Assumptions 1, 2 and 4a). While Assumptions 1, 2 can be verified similarly as in Example 1, the verification of Assumption 4a) is similar to [54].

TABLE 3

Example 2: Percentage of correct number of estimated threshold, and correct model structure specification. MDL, minimum description length; BIC, Bayesian information criterion.

Size	Regime structure classified(%), MDL	Model structure classified(%), MDL	Regime structure classified(%), BIC	Model structure classified(%), BIC
200	83.3	5.00	56.0	3.00
500	97.3	47.7	78.0	36.7
1000	99.3	79.3	97.7	78.7
2000	100.0	90.0	100.0	88.0

TABLE 4

Example 2: Mean and standard derivations of parameter estimates, thresholds and autoregressive coefficients. Parentheses: empirical standard derivations. The threshold estimates are based on replications with correct number of regimes. The coefficient estimates are based on replications with correct model structure.

Size	$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\omega}_1$	$\hat{\psi}_{y,1,1}^{(1)}$	$\hat{\omega}_2$	$\hat{\psi}_{y,1,2}^{(1)}$	$\hat{\omega}_3$	$\hat{\psi}_{y,1,3}^{(1)}$	$\hat{\psi}_{y,1,3}^{(2)}$
200	-0.213 (0.122)	0.402 (0.063)	-0.143 (0.057)	0.565 (0.096)	0.107 (0.028)	0.369 (0.086)	0.244 (0.177)	0.196 (0.110)	-0.708 (0.093)
500	-0.207 (0.030)	0.402 (0.018)	-0.154 (0.036)	0.594 (0.067)	0.101 (0.018)	0.267 (0.061)	0.294 (0.096)	0.253 (0.078)	-0.696 (0.079)
1000	-0.202 (0.016)	0.400 (0.008)	-0.148 (0.026)	0.599 (0.049)	0.101 (0.009)	0.250 (0.047)	0.300 (0.064)	0.250 (0.055)	-0.703 (0.061)
2000	-0.201 (0.008)	0.399 (0.004)	-0.149 (0.020)	0.600 (0.037)	0.100 (0.007)	0.248 (0.033)	0.297 (0.040)	0.247 (0.040)	-0.697 (0.040)
True	-0.2	0.4	-0.15	0.6	0.1	0.25	0.3	0.25	-0.7

TABLE 5

Example 2: Mean and standard derivations of parameter estimates, log-link coefficients. Parentheses: empirical standard derivations. The coefficient estimates are based on replications with correct model structure.

Size	$\hat{\alpha}_1$	$\hat{\psi}_{y,2,1}^{(1)}$	$\hat{\alpha}_2$	$\hat{\psi}_{y,2,2}^{(1)}$	$\hat{\psi}_{y,2,2}^{(2)}$	$\hat{\alpha}_3$
200	-0.401(0.176)	0.359(0.067)	-0.099(0.242)	0.447(0.077)	0.333(0.096)	0.016(0.138)
500	-0.419(0.135)	0.357(0.044)	-0.213(0.168)	0.455(0.047)	0.253(0.049)	0.114(0.112)
1000	-0.425(0.093)	0.349(0.033)	-0.205(0.123)	0.451(0.032)	0.254(0.033)	0.139(0.074)
2000	-0.404(0.067)	0.350(0.025)	-0.202(0.092)	0.449(0.022)	0.251(0.022)	0.145(0.058)
True	-0.4	0.35	-0.2	0.45	0.25	0.15

The empirical classification rate of regime and model structure for (6.2), and comparisons between criteria, are reported in Table 3. Again, the results are promising for moderately large sample sizes, and minimal description length still achieves superior performance. Tables 4 and 5 summarize estimation results for thresholds and model parameters, where the asymptotic convergences such

as  $O_p(n^{-1})$  rate for  $\hat{\Theta}$ , and  $O_p(n^{-1/2})$  for  $\hat{\Psi}$ , are recognized. In conclusion, the purposed methodology has satisfactory performance.

## 7. Application

Initial public offerings (IPOs) are one of the most important funding sources in finance. IPO activities are found to be time-varying; see, for example, [50], which discusses the cyclical effect with respect to stock market bull/bear trends and cross-year variation of levels of monthly IPO volumes. This is evidenced by a large variation with clustering of small and large observations across yearly periods. For instance, retreats of IPO activities are observed in 1973-1979 for the oil crisis, in 1982 for the energy crisis, in 2001 for the burst of Internet bubbles, and in 2008-2009 for the financial crisis; meanwhile, fervent IPO activity occurred in 1983-1987 during the stable global market with economic expansion, and in 1992-2000 during the high-tech boom.

However, few studies of IPO activities have considered quantitative modeling. [39] and [40] purposed autoregressive models which incorporates past monthly initial returns and market participation proxies. Nevertheless, linear autoregressive modeling is theoretically questionable for integer IPO volumes. Moreover, as indicated in [49] and [68], different market scenarios exist in the IPO market. These facts suggest the necessity of regime classifications for proper modeling.

Thus, for theoretical soundness and modeling flexibility, the GTLVM is applied to IPO volumes modeling. We model the U.S. monthly net IPO volumes from January 1976 to March 2014, where the dataset is available at [51]) ([https://www.quandl.com/data/RITTER/US\\_IPO\\_STATS-Historical-US-IPO-Statistics](https://www.quandl.com/data/RITTER/US_IPO_STATS-Historical-US-IPO-Statistics)). By the definition in [51], net IPO volumes exclude issuance of penny stocks, units and close-end funds.

To flexibly model counting data, we use a negative binomial distribution with a canonical log-link function in order to capture different dispersions across regimes. Here

$$f_i(y_t; \lambda_t) = \sum_{i=1}^{r+1} \left[ \frac{\Gamma(y_t + k_i)}{\Gamma(k_i)\Gamma(y_t + 1)} \left( \frac{k_i}{k_i + \lambda_t} \right)^{k_i} \left( \frac{\lambda_t}{k_i + \lambda_t} \right)^{y_t} \right] I(y_{t-d} \in (\theta_{i-1}, \theta_i]),$$

$$y_t \in \mathbb{N},$$

where  $\lambda_t$  satisfies  $\lambda_t = E(y_t)$ , and  $k_i$  is the dispersion parameter. With  $y_{t-j}^* = \log[\max(y_{t-j}, 0.01)] = g_{y,2}(y_{t-j})$ , the link function is expressed as

$$\log(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,i}^{(j)} y_{t-j}^* + \sum_{k=1}^{p_X} \psi_{x,i}^{(k)} x_k \right] I(y_{t-d} \in (\theta_{i-1}, \theta_i]).$$

To capture the dependence of the variables in the previous 12 months, we set the maximum delay as  $D = 12$  and autoregressive order as  $Q_{Y,2} = 12$ . And as indicated by [39], IPO activities are affected by new information arriving during

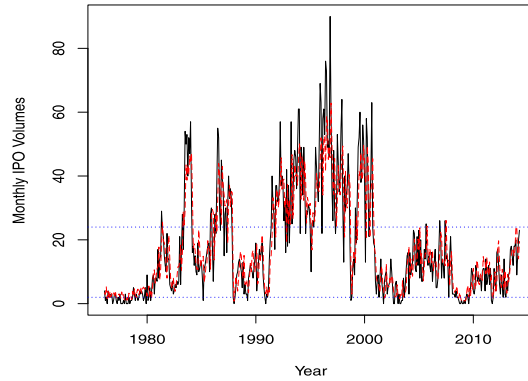


FIG 2. Plot of U.S. monthly net IPO volumes. Solid: original series; dashed: fitted series; dotted: threshold estimates

the book-building period which lasts for approximately two to four months. Hence, we check the averaged historical two-, three- or four-month return of the S&P 500 Index, denoted as  $x^{(2)}$ ,  $x^{(3)}$ , and  $x^{(4)}$ , respectively, as possible covariates that represent recent market performance. In addition, the past observations  $y_{t-2}^*$ ,  $y_{t-3}^*$  and  $y_{t-4}^*$  are included in the model. In each regime, the parameter  $\Psi_i$  is estimated by quasi-maximum likelihood with iteratively reweighted least square method, see [42].

For the net IPO series, two thresholds are estimated as  $\hat{\theta}_1 = 2$  and  $\hat{\theta}_2 = 24$ . The estimated dispersion  $k_i$  in the three regimes are 1.5962, 4.1817 and 9.9228, with theoretical standard errors of 0.4593, 0.5342 and 1.6211, respectively. The link function estimate is

$$\log(\lambda_t) = \begin{cases} 1.0242 + 0.1487y_{t-4}^* + 11.9900x_t^{(2)}, & y_{t-1} \leq 2, \\ (0.1178) (0.0484) & (2.9743) \\ 0.4552 + 0.5595y_{t-1}^* + 0.1549y_{t-2}^* + 0.1300y_{t-4}^* + 7.6926x_t^{(2)}, & \\ (0.1508) (0.0746) & (0.0360) & (0.0290) & (1.3092) \\ 2 < y_{t-1} \leq 24, & \\ 1.0459 + 0.5608y_{t-1}^* + 0.1468y_{t-4}^*, & y_{t-1} > 24. \\ (0.4239) (0.1222) & (0.0705) \end{cases} \tag{7.1}$$

The plot of fitted values is displayed in Figure 2. For model diagnostics, the standardized deviance residuals are plotted in Figure 3. The fluctuations of deviance residuals around zero indicates that the fitting is adequate.

Some discussions about the estimation results are as follows. First, [49] asserts that IPO activities can be classified in to two regimes: “cold” and “hot” markets, depending on the volumes of the activities. Later, [68] propose a more sophisticated three-regimes classification in terms of “cold”, “normal” and “hot”

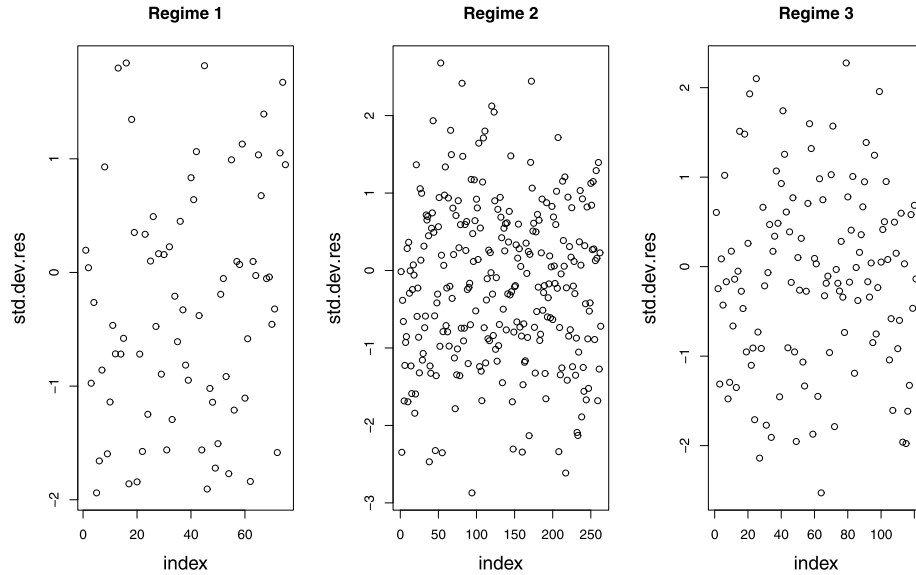


FIG 3. Model diagnostic plot: standardized deviance residuals of net IPO series for the three estimated regimes.

markets. The results in (7.1) suggest that the classification in [68] is more appropriate.

Second, (7.1) indicates that stock market return is an effective predictor of IPO volumes. In particular, the positive coefficients of  $x^{(2)}$  in regimes 1 and 2 of (7.1) indicate that IPO activities are positively associated with stock market performances. This phenomenon agrees with the theory in [45] that high market returns increase the incentives of IPO issuance, thus contributing to IPO market activity as volumes soar. Moreover, the coefficient of  $x^{(2)}$  is decreasing from regime 2 to regime 1, and becomes insignificant in regime 3, indicating that the positive effect of stock market returns diminishes with the increase in recent IPO activities. One possible explanation for this is as follows: as mentioned in [49], when the market is overactive, severe underpricing exists and discourages entrepreneurs. Hence, entrepreneurs choose to issue stocks in other periods, which offsets the market performance influence.

### Acknowledgment

We would like to thank the Editor, an Associate Editor and an anonymous reviewer for their helpful comments and useful suggestions, which greatly improve the presentation of this paper. Li's research was supported by the Fundamental Research Funds for the Central Universities in UIBE (CXTD10-09). Yau's research was supported in part by grants from HKSAR-RGC-GRF Nos 405113, 14305517 and 14601015.

## Appendix

### Proof of Theorem 3.1 (strict stationarity and ergodicity)

First, we state the following definitions from Markov chain theory that provide the background for studying strict stationarity and ergodicity.

**Definition 1. Irreducibility** A Markov process  $\{Y_t\}$  on a measurable space  $\{\Omega, \mathcal{B}\}$  with transition probability  $P^n(Y, A) = \text{pr}(Y_n \in A \mid Y_0 = y)$  is said to be  $\mu$ -irreducible for a measure  $\mu$  on  $\mathcal{B}$  if  $\sum_{n=1}^\infty P^n(y, A) > 0$  for all  $y \in \Omega$  whenever  $\mu(A) > 0$ .

**Definition 2. Geometric ergodicity:** a Markov process  $\{Y_t\}$  on  $\{\Omega, \mathcal{B}\}$  is geometrically ergodic if there exists a probability measure  $\pi$  on  $\mathcal{B}$  such that for all  $A \in \mathcal{B}$  and  $y \in \Omega$ , there exist  $\rho \in (0, 1)$  and  $M_y > 0$  that

$$\|P^n(y, A) - \pi(A)\| \leq \rho^n M_y,$$

where  $\|\cdot\|$  is the total-variation norm. This implies that  $\{y_t\}$  is ergodic,  $\beta$ -mixing and has a unique stationary distribution  $\pi$ ; see [4] and [52].

**Definition 3. Small set and petite set:** a set  $C \in \mathcal{B}$  is said to be small if there exists an integer  $m > 0$  and a non-trivial measure  $v_m(\cdot)$  on  $\mathcal{B}$  such that for all  $y \in C$  and  $A \in \mathcal{B}$ ,  $P^m(y, A) \geq v_m(A)$ .

Similarly, a set  $C$  is said to be petite for  $\{y_t\}$  if there exists a probability measure  $\gamma^*(\cdot)$  on  $\mathbb{N}^+$  and a non-trivial measure  $v_{\gamma^*}(\cdot)$  on  $\mathcal{B}$  such that for all  $y \in C$  and  $A \in \mathcal{B}$ ,  $\sum_{n=0}^\infty P^n(y, A)\gamma^*(n) \geq v_{\gamma^*}(A)$ . Clearly, a small set is a petite set; see [43].

The proof of Theorem 3.1 relies on the following theorem about the ergodicity of Markov chains.

**Theorem 7.1.** ([43]) *Let  $n(y) : \Omega \rightarrow \mathbb{N}^+$  be an integer valued function. An irreducible chain  $\{Y_t\}$  on  $\Omega$  is geometrically ergodic if it is aperiodic and there exists a non-negative function  $V \geq 1$  on  $\Omega$  which is bounded on a petite set  $C$ , and for all  $y \in \Omega$ , there exist  $\rho \in (0, 1)$  and  $b \in (0, \infty)$  such that*

$$\begin{aligned} \mathbb{E}[V(Y_{t+n(Y_t)}) \mid Y_t] &= \int P^{n(Y_t)}(Y_t, dY_{t+n(Y_t)})V(Y_{t+n(Y_t)}) \\ &\leq \rho^{n(Y_t)}[V(Y_t) + bI(Y_t \in C)]. \end{aligned} \tag{7.2}$$

By Condition 1a), denote

$$\begin{aligned} p^* &= \max\{p_{Y,1}, p_{Y,2}, \tilde{p}, p_e, p_\epsilon, q + d, 1\}, \quad \xi_t = (y_t, X_t, e_t, \epsilon_t), \\ \Xi_t &= \{(\xi_t, \dots, \xi_{t-p^*+1})\}, \end{aligned}$$

and  $\mathcal{F}_{t-1}$  as the sigma-field generated by  $\{\Xi_{t-1}, \Xi_{t-2}, \dots\}$ . First, we show that  $\{\Xi_t\}$  is Markovian under Condition 1a). Note that

$$\begin{aligned} \text{pr}(\{X_t, \dots, X_{t-p^*+1}\} \mid \mathcal{F}_{t-1}) &= \text{pr}(\{X_t, \dots, X_{t-p^*+1}\} \mid \{X_{t-1}, \dots, X_{t-p^*}\}) \\ &= \text{pr}(\{X_t, \dots, X_{t-p^*+1}\} \mid \Xi_{t-1}). \end{aligned} \tag{7.3}$$

Also, as  $z_{t-d}$  is measurable with respect to the sigma-field generated by  $\{X_{t-d}, y_{t-d}, \dots, X_{t-d-q}, y_{t-d-q}\}$ , we have that  $\text{pr}(z_{t-d} \mid \mathcal{F}_{t-1}) = \text{pr}(z_{t-d} \mid \Xi_{t-1})$ . Hence,

$$\begin{aligned} \text{pr}(\Xi_t, z_{t-d} \mid \mathcal{F}_{t-1}) &= \text{pr}(z_{t-d} \mid \mathcal{F}_{t-1})\text{pr}(\Xi_t \mid \mathcal{F}_{t-1}, z_{t-d}) \\ &= \text{pr}(z_{t-d} \mid \Xi_{t-1})\text{pr}(\Xi_t \mid \Xi_{t-1}, z_{t-d}) = \text{pr}(\Xi_t, z_{t-d} \mid \Xi_{t-1}). \end{aligned}$$

Integrating both sides with respect to the density of  $z_{t-d}$ , we have  $\text{pr}(\Xi_t \mid \mathcal{F}_{t-1}) = \text{pr}(\Xi_t \mid \Xi_{t-1})$ , and hence  $\{\Xi_t\}$  is a Markov process.

Under Condition 1a'), we can analogously verify the Markovian property of  $\{\Xi_t, z_t^*\}$ . Denote  $\mathcal{F}_{t-1}$  as the sigma-field generated by  $\{\Xi_{t-1}, z_{t-1}^*, \Xi_{t-2}, z_{t-2}^*, \dots\}$ . As (7.3) and  $\text{pr}(z_t^* \mid \mathcal{F}_{t-1}) = \text{pr}(z_t^* \mid \Xi_{t-1}, z_{t-1}^*)$  hold by Condition 1a'), we have

$$\begin{aligned} \text{pr}(\Xi_t, z_t^* \mid \mathcal{F}_{t-1}) &= \text{pr}(z_t^* \mid \mathcal{F}_{t-1})\text{pr}(\Xi_t \mid \mathcal{F}_{t-1}, z_t^*) \\ &= \text{pr}(z_t^* \mid \Xi_{t-1}, z_{t-1}^*)\text{pr}(\Xi_t \mid \Xi_{t-1}, z_t^*) \\ &= \text{pr}(z_t^* \mid \Xi_{t-1}, z_{t-1}^*)\text{pr}(\Xi_t \mid \Xi_{t-1}, z_t^*, z_{t-1}^*) \\ &= \text{pr}(\Xi_t, z_t^* \mid \Xi_{t-1}, z_{t-1}^*). \end{aligned}$$

Thus,  $\{\Xi_t, z_t^*\}$  is shown to be Markovian. Next, we illustrate our proof with  $p_X = p_e = p_\epsilon = 0$  for simplicity. In this case, it suffices to prove stationarity and ergodicity of  $\{Y_t\} = (y_t, \dots, y_{t-p^*+1})$ . We will show the geometric ergodicity of  $\{y_t\}$  by using Theorem 7.1. Hence, we need to verify that  $\{Y_t\}$  or  $\{Y_t, Z_t^*\}$  is an irreducible and aperiodic Markov process and construct the corresponding function  $V$ . For the general case (except for the irreducibility when  $p_e > 0$ ), the same verification methodologies could be applied on  $\{\Xi_t\}$  or  $\{\Xi_t, z_t^*\}$  with mild modifications, and hence the proof is omitted. Therefore, it suffices to show the stationarity and ergodicity of

$$\begin{cases} y_t = \sum_{i=1}^{r+1} \left[ \omega_i + \sum_{l=1}^{p_{Y,1}} \psi_{y,1,i}^{(l)} g_{y,1}(y_{t-l}) + g_\lambda(\lambda_t, u_t) + \psi_{\epsilon,i}^{(0)} \epsilon_t \right] \\ \quad \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]), \\ \phi(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y_{t-j}) + \psi_{\epsilon,i}^{(0)} \epsilon_t \right] I(z_{t-d} \in (\theta_{i-1}, \theta_i]). \end{cases} \quad (7.4)$$

Denote  $\mu$  as the Lebesgue measure on  $\mathbb{R}$ , and  $\mu^{p^*}$  as the Lebesgue measure on  $\mathbb{R}^{p^*}$ . Next, we show that  $\{Y_t\}$  or  $\{Y_t, Z_t^*\}$  is irreducible and aperiodic, and there exists some small set by the following proposition.

**Proposition 1.** *Under (7.4), we have:*

1. *Under Condition 1a),  $\{Y_t\}$  is  $\mu^{p^*}$ -irreducible and aperiodic. In addition, sets of the form  $C = \{Y_t : |Y_t|_\infty \leq c\}$  for some  $c > 0$  is small for  $\{Y_t\}$ .*
2. *Under Condition 1a'),  $\{Y_t, Z_t^*\}$  is  $\mu^{p^*} \times \nu^q$ -irreducible with some discrete measure  $\nu$  on  $\{0, 1\}^r$  for  $r$ -manifolds of set  $\{0, 1\}$ , and  $\{Y_t, Z_t^*\}$  is aperiodic. In addition, sets of the form  $C = \{(Y_t, Z_t^*) : |Y_t|_\infty \leq c, |Z_t^*|_\infty \leq q\}$  for some  $c > 0$  is small for  $\{Y_t, Z_t^*\}$ .*

*Proof of Proposition 1.* We illustrate the proof of  $\{Y_t\}$  under Condition 1a), where the proof of  $\{Y_t, Z_t^*\}$  under condition 1a') follows similar arguments. First we assume that  $e_t$  and  $\epsilon_t$  have almost everywhere continuous and positive densities on  $\mathbb{R}$ . We divide the proof into three parts: irreducibility, existence of small set, and aperiodicity.

1) *Irreducibility.* As  $e_t$  has positive density on  $\mathbb{R}$ ,  $y_t$  can reach any point on  $\mathbb{R}$ . Let  $A \subset \mathbb{R}$  satisfies  $0 < \mu(A) < \infty$ , we first show that  $\text{pr}(y_t \in A \mid Y_{t-1}, \lambda_t, u_t, z_{t-d})$  is positive for any realization  $\{Y_{t-1}, \lambda_t, u_t, z_{t-d}\} = \{(y_{t-1}, \dots, y_{t-p^*}), \lambda_t, u_t, z_{t-d}\}$  with  $z_{t-d} \in (\theta_{i-1}, \theta_i]$ . With respect to (7.4), we can construct

$$E_{t,i} = \left\{ e_t : \left[ \alpha_i + \sum_{l=1}^{p_{Y,1}} \psi_{y,1,i}^{(l)} g_{y,1}(y_{t-l}) + g_\lambda(\lambda_t, u_t) + \psi_{e,i}^{(0)} e_t \right] \in A \mid Y_{t-1}, \lambda_t, u_t, z_{t-d} \right\}.$$

Hence, an injection exists between  $E_{t,i}$  and  $A$  and thus  $\mu(E_{t,i}) > 0$ . The mapping from  $E_{t,i}$  to  $A$  is surjective and  $\mu(E_{t,i}) > 0$  holds. Denote the density of  $e_t$  as  $f_e(\cdot)$ , we have

$$\begin{aligned} \text{pr}(y_t \in A \mid Y_{t-1}, \lambda_t, u_t, z_{t-d}) &= \text{pr}(e_t \in E_{t,i} \mid Y_{t-1}, \lambda_t, u_t, z_{t-d}) \\ &= \int_{v \in E_{t,i}} f_e(v) dv \geq \inf_{e_t \in E_{t,i}} f_e(e_t) \mu(E_{t,i}) > 0. \end{aligned}$$

As  $e_t$  and  $\epsilon_t$  have almost everywhere positive densities, the conditional densities  $f_i(y_t \mid Y_{t-1}, \lambda_t, u_t, z_{t-d})$  and  $h_i(\lambda_t \mid Y_{t-1}, z_{t-d})$  are thus almost everywhere positive. Therefore, from (2.1), the marginally density  $f$  is positive almost everywhere. Denote the marginal density of  $u_t$  and  $\epsilon_t$  as  $f_u(\cdot)$  and  $f_\epsilon(\cdot)$ , respectively, it follows that

$$\begin{aligned} \text{pr}(y_t \in A \mid Y_{t-1}, z_{t-d}) &= \int \text{pr}(y_t \in A \mid Y_{t-1}, \lambda_t, z_{t-d}) h_i(\lambda_t \mid Y_{t-1}, z_{t-d}) d\lambda_t \\ &= \int \int_0^1 \text{pr}(y_t \in A \mid Y_{t-1}, \lambda_t, u, z_{t-d}) f_u(u) du h_i(\lambda_t \mid Y_{t-1}, z_{t-d}) d\lambda_t \\ &> 0, \end{aligned} \tag{7.5}$$

for some set  $B$  with  $\mu(B) > 0$ , where  $f_u(u) = 1$  on  $[0, 1]$  and  $z_{t-d} \in (\theta_{i-1}, \theta_i]$ . Furthermore, denote  $f_z(\cdot)$  as the density of  $z_t$ , we have, almost surely,

$$\text{pr}(y_t \in A \mid Y_{t-1}) = \int \text{pr}(y_t \in A \mid Y_{t-1}, z_{t-d}) f_z(z_{t-d} \mid Y_{t-1}) dz_{t-d} > 0. \tag{7.6}$$

Therefore, for any  $\tilde{A} \subset \mathbb{R}^{p^*}$  and  $y \in \mathbb{R}^{p^*}$  with Lebesgue measure  $\mu_{p^*}(\tilde{A}) > 0$ , as there exist some  $A_1 \in \mathbb{R}$  and  $A_2 \in \mathbb{R}^{p^*-1}$  such that  $A_1 \times A_2 \subset \tilde{A}$  and  $\mu(A_1) > 0$ ,

we have  $P(Y_t \in \tilde{A} | Y_{t-1} = y) > P(y_t \in A_1 | Y_{t-1} = y) > 0$ . This completes the proof of irreducibility.

2) *Existence of a small set  $C$ .* By Definition 3, we need to construct a small set  $C$  such that, for all  $Y_0 = y \in C \subset \mathbb{R}^{p^*}$  with  $\mu_{p^*}(C) > 0$ , there exists a non-trivial measure  $v_m(\cdot)$  on the Borel sigma-field on  $\mathbb{R}^{p^*}$ ,  $\mathcal{B}$ , such that for any  $\tilde{A} \in \mathcal{B}$  with  $\mu_{p^*}(\tilde{A}) > 0$ ,

$$P^m(y, \tilde{A}) \geq v_m(\tilde{A}) > 0.$$

We construct the set as  $C = \{y : y = (y_1, \dots, y_{p^*}) \in \mathbb{R}^{p^*}, |y|_\infty \leq c\}$  for some constant  $c$ , where  $|y|_\infty = \max\{|y_1|, \dots, |y_{p^*}|\}$ . First, as the density of  $e_t$  is almost everywhere positive,  $f(y_t | Y_{t-1}, \lambda_t, u_t, z_{t-d}) > 0$ . Analogous to (7.5)–(7.6), we have that  $f(y_t | Y_{t-1}) > 0$ , and hence with (7.6),

$$\begin{aligned} \text{pr}(y_{t+1} \in A | Y_{t-1}) &= \int \text{pr}(y_{t+1} \in A | Y_{t-1}, y_t) f(y_t | Y_{t-1}) dy_t \\ &= \int \text{pr}(y_{t+1} \in A | Y_t) f(y_t | Y_{t-1}) dy_t > 0, \end{aligned} \quad (7.7)$$

for any  $A \subset \mathbb{R}$  with  $\mu(A) > 0$ . By using induction on (7.7),  $\text{pr}(y_{t+m-1} \in A | Y_{t-1}) > 0$  and hence  $\text{pr}(Y_{t+m-1} \in \tilde{A} | Y_{t-1}) > 0$  for all positive integer  $m$  and any  $\tilde{A} \in \mathcal{B}$  with  $\mu_{p^*}(\tilde{A}) > 0$ . Thus, setting

$$v_m(\tilde{A}) = \min_{y \in C} \text{pr}(y_{m-1} \in \tilde{A} | Y_0 = y) = \min_{y \in C} P^m(y, \tilde{A}),$$

we have  $v_m(\tilde{A}) > 0$  by the compactness of  $C$ . Thus, the set  $C$  is verified as a small set.

3) *Aperiodicity.* From the existence proof of the small set  $C$ , it has been shown that  $P^1(y, C) > 0$  and  $P^2(y, C) > 0$  for all  $y \in C$ . By Proposition A1.1 of [9], it follows that  $\{y_t\}$  is aperiodic.

To relax the assumption that  $\epsilon_t$  and  $e_t$  have almost everywhere continuous and positive densities on  $\mathbb{R}$ , we extend to the case  $e_t = \epsilon_t = 0$  for all  $t$  where  $e_t$  and  $\epsilon_t$  do not have almost everywhere positive density. Define a perturbation  $\{y_t^m\}$  of  $\{y_t\}$  by

$$\begin{cases} y_t^m = \sum_{i=1}^{r+1} \left[ \omega_i + \sum_{l=1}^{p_{Y,1}} \psi_{y,1,i}^{(l)} g_{y,1}(y_{t-l}^m) + g_\lambda(\lambda_t, u_t) \right] I(z_{t-d} \in (\theta_{i-1}, \theta_i]) \\ \quad + \sigma_{1,m} e_{1,t}, \\ \phi(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y_{t-j}^m) \right] I(z_{t-d} \in (\theta_{i-1}, \theta_i]) + \sigma_{2,m} e_{2,t}, \end{cases} \quad (7.8)$$

where  $\sigma_{1,m}, \sigma_{2,m} > 0, \sigma_{1,m} \rightarrow 0$  and  $\sigma_{2,m} \rightarrow 0$  as  $m \rightarrow \infty$ , and  $e_{1,t}, e_{2,t}$  are i.i.d. zero-mean noises with finite first moment and almost everywhere positive densities. Using the perturbation techniques in [22],  $\{y_t^m\}$  is irreducible, aperiodic, and a small set exists. Thus we can derive the strict stationarity and ergodicity of  $\{y_t^m\}$  by Theorem 7.1. Since  $\{y_t^m\}$  converges almost surely to  $\{y_t\}$  as  $m \rightarrow \infty$ , the proof is complete.  $\square$



*Proof of Lemma 3.1.* Lemma 3.1 is established in the irreducibility part of the proof of Proposition 1.  $\square$

*Proof of Theorem 3.1.* First, we give the proof with respect to  $\{Y_t\}$  under Condition 1a). Recall that the irreducibility, aperiodicity, and existence of a small set  $C = \{y : y = (y_1, \dots, y_{p^*}) \in \mathbb{R}^{p^*}, |y|_\infty \leq c\}$  have been established in Proposition 1. From Theorem 7.1, it suffices to construct a function  $V$  that satisfies (7.2). For the model (7.4), let  $Y_{t-1} = (y_{t-1}, \dots, y_{t-p^*})$  and define

$$V(Y_{t-1}) = 1 + |Y_{t-1}|_\infty, \tag{7.9}$$

where  $|Y_{t-1}|_\infty = \max\{|y_{t-1}|, \dots, |y_{t-p^*}|\}$ . Conditional on  $Y_{t-1}$ , we have

$$\begin{aligned} \mathbb{E}(|y_t| \mid Y_{t-1}) &\leq \mathbb{E} \left[ \left\{ \sum_{i=1}^{r+1} |\omega_i| + \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| |g_{y,1}(y_{t-l})| + |g_\lambda(\lambda_t, u_t)| + |\psi_{\epsilon,i}^{(0)} e_t| \right\} \right. \\ &\quad \left. \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]) \mid Y_{t-1} \right] \\ &\leq \max_i \left\{ \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| |g_{y,1}(y_{t-l})| + \mathbb{E}(|g_\lambda(\lambda_t, u_t)| \mid Y_{t-1}) \right\} \\ &\quad + \max_i (|\omega_i| + \mathbb{E}|\psi_{\epsilon,i}^{(0)} e_t|). \end{aligned} \tag{7.10}$$

Denote  $y^* = \arg \max_{j=1 \dots p_{Y,2}} g_{y,2}(y_{t-j})$ . For the case of concave  $\phi(\cdot)$ , Condition 1b) implies that the inverse function  $\phi^{-1}(\cdot)$  exists, and is strictly increasing and convex. Hence, as  $\max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| < 1$ ,

$$\begin{aligned} \lambda_t &= \phi^{-1} \left\{ \sum_{i=1}^{r+1} \left( \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y_{t-j}) + \psi_{\epsilon,i}^{(0)} \epsilon_t \right) I(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right\} \\ &\leq \phi^{-1} \left\{ \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| g_{y,2}(y^*) + \max_i (\alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t|) \right\} \\ &\leq \left( 1 - \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| \right) \phi^{-1} \left\{ \frac{\max_i (\alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t|)}{1 - \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}|} \right\} \\ &\quad + \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| \phi^{-1}(g_{y,2}(y^*)) \\ &\leq \left( 1 - \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| \right) \phi^{-1} \left\{ \frac{\max_i (\alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t|)}{1 - \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}|} \right\} \\ &\quad + \max_i \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| \max\{b_2 |Y_{t-1}|_\infty, G_2\}, \end{aligned}$$

where the second last inequality follows from Jensen’s inequality, and the last inequality follows from Condition 1c). Let  $\mathcal{F}_{t-1}$  be the  $\sigma$ -field generated by  $\{y_s, \epsilon_s\}_{s \leq t-1}$ . It follows from  $E\{\phi^{-1}(\epsilon_t)\} < \infty$  that

$$\begin{aligned}
 & E(g_\lambda(\lambda_t, u_t) \mid Y_{t-1}) \\
 & \leq E \left[ \left| g_\lambda \left\{ \phi^{-1} \left( \max_i \left( \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y^*) + \psi_{\epsilon,i}^{(0)} \epsilon_t \right) \right), u_t \right\} \right| \mid Y_{t-1} \right] \\
 & = E \left[ E \left[ \left| g_\lambda \left\{ \phi^{-1} \left( \max_i \left( \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y^*) + \psi_{\epsilon,i}^{(0)} \epsilon_t \right) \right) \right\} \right| \right. \right. \\
 & \quad \left. \left. \mid \mathcal{F}_{t-1} \right] \mid Y_{t-1} \right] \\
 & \leq E \left[ \max \left\{ H, b_\lambda \left( 1 - \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right| \right) \phi^{-1} \left( \frac{\max_i \left( \alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t| \right)}{1 - \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right|} \right) \right. \right. \\
 & \quad \left. \left. + b_\lambda \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right| \max\{b_2 |Y_{t-1}|_\infty, G_2\} \right\} \mid Y_{t-1} \right] \\
 & \leq b_\lambda b_2 \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right| |Y_{t-1}|_\infty + \tilde{H}^*, \tag{7.11}
 \end{aligned}$$

for some constant  $\tilde{H}^*$ . Equipped with (7.11), it can be shown from (7.10) that

$$\begin{aligned}
 E(|y_t| \mid Y_{t-1}) & \leq \max_i \left( \sum_{l=1}^{p_{Y,1}} \left| \psi_{y,1,i}^{(l)} \right| |g_{y,1}(y_{t-1})| + b_\lambda b_2 \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right| |Y_{t-1}|_\infty \right) \\
 & \quad + \tilde{H}^* + \max_i \left( |\omega_i| + E \left| \psi_{\epsilon,i}^{(0)} \epsilon_t \right| \right) \\
 & \leq H^* + \max_i \rho_i(1) |Y_{t-1}|_\infty = H^* + \tilde{\rho} |Y_{t-1}|_\infty, \tag{7.12}
 \end{aligned}$$

where  $\tilde{\rho} = \max_i \rho_i(1) \in (0, 1)$  and  $H^*$  is some constant. Using (7.12), we have

$$\begin{aligned}
 E(|y_{t+1}| \mid Y_{t-1}) & = E\{E(|y_{t+1}| \mid y_t, Y_{t-1}) \mid Y_{t-1}\} \\
 & \leq H^* + \tilde{\rho} \max\{|Y_{t-1}|_\infty, E(|y_t| \mid Y_{t-1})\} \\
 & \leq (1 + \tilde{\rho})H^* + \tilde{\rho} |Y_{t-1}|_\infty.
 \end{aligned}$$

Arguing inductively, for  $Y_{t+p^*-1} = (y_{t+p^*-1}, \dots, y_t)$ , we have

$$E(|Y_{t+p^*-1}|_\infty \mid Y_{t-1}) \leq (1 + \tilde{\rho} + \dots + \tilde{\rho}^{p^*-1})H^* + \tilde{\rho} |Y_{t-1}|_\infty. \tag{7.13}$$

Select  $\rho \in (\tilde{\rho}^{1/p^*}, 1)$ . Then, taking  $n(Y_{t-1}) = p^*$ ,  $b = [(1 + \tilde{\rho} + \dots + \tilde{\rho}^{p^*-1})H^* + 1]/\rho^{p^*} - 1$  and  $c = \max\{\tilde{y}, b/(\rho^{p^*} - \tilde{\rho})\}$  for the small set  $C$ , it can be verified that

$V$  in (7.9) satisfies (7.2). By Theorem 7.1,  $\{y_t\}$  is geometric ergodicity. Thus, by Definition 2,  $\{y_t\}$  is strict stationary and ergodic.

For the case that  $\phi(\cdot)$  is a polynomial with order  $\gamma$ , we have  $\phi(\lambda) = \lambda^\gamma(1 + o(1))$  and so  $\phi^{-1}(v) = v^{1/\gamma}(1 + o(1))$ . For  $\gamma \geq 1$ ,

$$\begin{aligned} \lambda_t &= \phi^{-1} \left\{ \sum_{i=1}^{r+1} \left( \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y_{t-j}) + \psi_{\epsilon,i}^{(0)} \epsilon_t \right) I(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right\} \\ &\leq \left\{ \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right| g_{y,2}(y^*) + \max_i \left( \alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t| \right) \right\}^{1/\gamma} (1 + o(1)) \\ &\leq \left\{ \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right|^{1/\gamma} (g_{y,2}(y^*))^{1/\gamma} + \max_i \left( \alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t| \right)^{1/\gamma} \right\} (1 + o(1)) \\ &\leq \left\{ \max_i \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right|^{1/\gamma} \max\{b_2 |Y_{t-1}|_\infty, G_2\} + \phi^{-1} \left( \max_i \left( \alpha_i + \psi_{\epsilon,i}^{(0)} |\epsilon_t| \right) \right) \right\} \\ &\qquad \qquad \qquad \times (1 + o(1)), \quad (7.14) \end{aligned}$$

where the last inequality follows from Condition 1c) and  $\phi^{-1}(v) = v^{1/\gamma}(1 + o(1))$ . The  $(1 + o(1))$  term in (7.14) is negligible if  $|Y_{t-1}|_\infty$  is large. Using the same arguments as in the derivation of (7.11), (7.12) and (7.13), the geometric ergodicity can be derived similarly.

The geometric ergodicity of  $\{Y_t, Z_t^*\}$  under condition 1a') can be shown analogously: By Proposition 1, sets of the form  $C = \{(Y_t, Z_t^*) : |Y_t|_\infty \leq c, |Z_t^*|_\infty \leq q\}$  for some  $c > 0$  are small for  $\{Y_t, Z_t^*\}$ . Define  $V(Y_t, Z_t^*) = 1 + |Y_t|_\infty + |Z_t^*|_\infty$ , we have that (7.13) holds by the same derivation. Then, since  $0 \leq |Z_t^*|_\infty \leq 1$ , by selecting  $b = [(1 + \tilde{\rho} + \dots + \tilde{\rho}^{p^* - 1})H^* + 1 + q]/\rho^{p^*} - 1$ ,  $c = \max\{\tilde{y}, b/(\rho^{p^*} - \tilde{\rho})\}$ , with other parameters the same as in the proof of  $\{Y_t\}$  under condition 1a), the geometric ergodicity of  $\{Y_t, Z_t^*\}$  is verified by Theorem 7.1. Hence,  $\{Y_t, Z_t^*\}$  is strictly stationary and ergodic, and so does  $y_t$ . This finishes the proof.  $\square$

*Proof of Corollary 1.* Denote

$$\rho_z(\gamma, z_{t-d}) = \sum_{i=1}^{r+1} \sum_{j=1}^{p_{Y,2}} \left| \psi_{y,2,i}^{(j)} \right|^{1/\gamma} I(z_{t-d} \in (\theta_{i-1}, \theta_i]).$$

Let  $\mathcal{F}_{t-1}$  be the  $\sigma$ -field generated by  $\{y_s, \epsilon_s\}_{s \leq t-1}$ . From the independence of  $z_{t-d}$  and  $\{y_s\}_{s < t}$ , for concave  $\phi$  in 1'), we have

$$\begin{aligned} &E(|g_\lambda(\lambda_t, u_t)| \mid Y_{t-1}) \\ &= E \left[ E \left[ \left| g_\lambda \left\{ \phi^{-1} \left( \sum_{i=1}^{r+1} \left( \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,2,i}^{(j)} g_{y,2}(y_{t-j}) + \psi_{\epsilon,i}^{(0)} \epsilon_{t-q} \right) \right) \right\} \right| \right] \right] \end{aligned}$$

$$\begin{aligned}
& \times I(z_{t-d} \in (\theta_{i-1}, \theta_i]) \Big) , u_t \Big| \Big| \mathcal{F}_{t-1} \Big| \Big| Y_{t-1} \Big] \\
& \leq \mathbb{E} \left[ \max \left\{ H, b_\lambda (1 - \rho_z(1, z_{t-d})) \phi^{-1} \left( \frac{\max_i (\alpha_i + \psi_{\epsilon, i}^{(0)} |\epsilon_{t-d}|)}{1 - \rho_z(1, z_{t-d})} \right) \right. \right. \\
& \qquad \qquad \qquad \left. \left. + b_\lambda \rho_z(1, z_{t-d}) \max\{b_2 |Y_{t-1}|_\infty, G_2\} \right\} \Big| Y_{t-1} \right] \\
& \leq b_\lambda b_2 \sum_{i=1}^{r+1} \left\{ \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| |Y_{t-1}|_\infty \Pr(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right\} + \tilde{H}_2^*, \tag{7.15}
\end{aligned}$$

for some constant  $\tilde{H}_2^*$ . Thus, it is straightforward to bound  $\mathbb{E}(|y_t| | Y_{t-1})$  by

$$\begin{aligned}
& \mathbb{E}(|y_t| | Y_{t-1}) \\
& \leq \sum_{i=1}^{r+1} \left[ \left\{ \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| |g_{y,1}(y_{t-l})| + b_\lambda b_2 \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| |Y_{t-1}|_\infty \right\} \Pr(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right] \\
& \quad + \tilde{H}_2^* + \max_i \left( |\omega_i| + \mathbb{E} |\psi_{e,i}^{(0)} e_{t-m}| \right) \\
& \leq H_2^* + \sum_{i=1}^{r+1} \left[ \left\{ b_1 \sum_{l=1}^{p_{Y,1}} |\psi_{y,1,i}^{(l)}| + b_\lambda b_2 \sum_{j=1}^{p_{Y,2}} |\psi_{y,2,i}^{(j)}| \right\} \Pr(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right] |Y_{t-1}|_\infty \\
& = H_2^* + \left\{ \sum_{i=1}^{r+1} \rho_i(1) \Pr(z_{t-d} \in (\theta_{i-1}, \theta_i]) \right\} |Y_{t-1}|_\infty \\
& \leq H_2^* + \tilde{\rho} |Y_{t-1}|_\infty,
\end{aligned}$$

for some  $\tilde{\rho} \in (0, 1)$  and some constants  $H_2^*$ . Using similar arguments in the proof of Theorem 3.1,  $\{y_t\}$  is geometrically ergodic and hence strictly stationary and ergodic. If  $\phi(\cdot)$  is a polynomial with order  $\gamma \geq 1$  in  $2'$ , then by considering  $\rho_z(\gamma, z_{t-d})$  instead of  $\rho_z(1, z_{t-d})$ , similar derivations yield an analogous result to (7.15), with  $|\psi_{y,2,i}^{(j)}|$  replaced by  $|\psi_{y,2,i}^{(j)}|^{1/\gamma}$ . Thus, the strict stationarity and ergodicity are verified.  $\square$

*Proof of Corollary 2.* First consider concave  $\phi$ . For simplicity, we first illustrate the proof for  $d = 1$  and assume  $v = \max_{i=2, \dots, r} \rho_i(1) > 1$ . Using similar arguments as in (7.12),  $\mathbb{E}(|y_t| | Y_{t-1}) \leq H^* + v |Y_{t-1}|_\infty$  for some  $H^*$ . Therefore,

$$\begin{aligned}
& \mathbb{E}(|y_{t+1}| | Y_{t-1}) \\
& = \mathbb{E}\{\mathbb{E}(|y_{t+1}| | y_t, Y_{t-1}) | Y_{t-1}\} \\
& \leq \mathbb{E} \left[ \sum_{i=1}^{r+1} \rho_i(1) I(y_t \in (\theta_{i-1}, \theta_i]) \max\{|Y_{t-1}|_\infty, |y_t|\} \Big| Y_{t-1} \right] + H^*
\end{aligned}$$

$$\begin{aligned}
 &\leq \left[ \max\{\rho_1(1), \rho_{r+1}(1)\} \Pr(\{y_t > \theta_r\} \cup \{y_t \leq \theta_1\} \mid Y_{t-1}) \right. \\
 &\quad \left. + \max_{i=2, \dots, r} \rho_i(1) \Pr(\theta_1 < y_t \leq \theta_r \mid Y_{t-1}) \right] \\
 &\qquad \qquad \qquad \times \max\{\mathbb{E}(|y_t| \mid Y_{t-1}), |Y_{t-1}|_\infty\} + H^* \\
 &\leq H^* + [\max\{\rho_1(1), \rho_{r+1}(1)\}(1 - \pi_y) + v\pi_y] (H^* + |vY_{t-1}|_\infty) \\
 &\leq (1 + \tilde{\rho})H^* + \tilde{\rho}v|Y_{t-1}|_\infty, \tag{7.16}
 \end{aligned}$$

for some  $\tilde{\rho} \in (0, 1)$ . Arguing inductively, analogous to (7.16), for all  $k_1 \geq 0$ ,  $0 \leq k_2 \leq p^* - 1$ ,

$$\mathbb{E}(|y_{t+k_1p^*+k_2+1}| \mid Y_{t-1}) \leq (1 + \tilde{\rho} + \dots + \tilde{\rho}^{k_1p^*+k_2+1})H^* + \tilde{\rho}^{k_1+1}v|Y_{t-1}|_\infty.$$

As  $\tilde{\rho}^{k_1+1}v < 1$  for sufficiently large  $k_1$ , the geometric ergodicity of  $\{y_t\}$  can be established using similar arguments in the proof in Theorem 3.1. When  $\phi$  is a polynomial of order  $\gamma \geq 1$ , we replace  $\rho_i(\gamma)$  by  $\rho_i(1)$  to derive the result analogous to (7.16). Using similar but more tedious arguments, the results can be derived for  $d = 2, \dots, p^*$  that,  $k_1 \geq 0$ ,  $0 \leq k_2 \leq p^* - 1$ ,

$$\begin{aligned}
 \mathbb{E}(|y_{t+k_1p^*+k_2+d}| \mid Y_{t-1}) &\leq (1 + \tilde{\rho} + \dots + \tilde{\rho}^{k_1p^*+k_2+1})(1 + v + \dots + v^{d-1})H^* \\
 &\quad + \tilde{\rho}^{k_1+1}v^d|Y_{t-1}|_\infty.
 \end{aligned}$$

Again, as  $\tilde{\rho}^{k_1+1}v^d < 1$  for sufficiently large  $k_1$ , the geometric ergodicity of  $\{y_t\}$  follows.  $\square$

*Proof of Corollary 3.* For the TARMA model, the threshold variable  $y_{t-d}$  is clearly measurable with respect to the sigma-field generated by  $\{X_t, y_t, \dots, X_{t-q}, y_{t-q}\}$  with  $q \geq d$ . Let  $b_1 = 1$ ,  $b_2 = 0$  and  $G_1 = \tilde{y} = 1$ . It follows that  $|g_{y,1}(y)| \leq b_1|y|$  for all  $|y| > \tilde{y}$ , and  $g_{y,1}(y) \leq G_1$  for all  $|y| \leq \tilde{y}$ . Together with the irreducibility from Condition 1f), Condition 1 holds. Therefore, by Theorem 3.1 and the fact that  $\psi_{y,2,i}^{(j)} = 0$  for all  $i$ , the TARMA process is stationary and ergodic if  $\max_{i=1, \dots, r+1} \rho_i(1) = \max_{i=1, \dots, r+1} \sum_{l=1}^{pY,1} |\psi_{y,1,i}^{(l)}| < 1$ . The same arguments apply to TAR model, with the irreducibility condition guaranteed by Lemma 3.1 instead of Condition 1f).  $\square$

**Verification of Assumption 4b)**

In this section we illustrate the verification of conditions in Assumption 4b) using the model in Example 1 of the simulation studies. First, we have  $T(y_t) = y_t$ . As Theorem 3.1 indicates that  $\{y_t\}$  is strictly stationary and ergodic,  $y_t < \infty$  almost surely. Denote  $Y_t = (y_t, \dots, y_{t-3})$  and let  $|\cdot|_\infty$  be the infinite norm such that  $|Y_t|_\infty = \max\{|y_t|, \dots, |y_{t-3}|\}$ . Note that  $\{Y_t\}$  is a Markov process. By (6.1),

$$\log(\lambda_t) \leq \max_i a_i + \max_i \sum_{j=1}^{p_i} |\psi_{i,j}| \log(|Y_{t-1}|_\infty + 1),$$

and hence

$$\lambda_t \leq e^{\max_i a_i} (|Y_{t-1}|_\infty + 1)^{\max_i \sum_{j=1}^{p_i} |\psi_{i,j}|} \leq K_\lambda |Y_{t-1}|_\infty^\rho + c_\lambda, \tag{7.17}$$

for some  $0 < \rho < 1$  and positive constants  $K_\lambda$  and  $c_\lambda$ . Thus,  $\lambda_t < \infty$  almost surely, and  $\text{var}(T(y_t) | \lambda_t) < \infty$  almost surely.

Next we show the existence of moments of  $y_t$  and  $\lambda_t$ . By strict stationarity and ergodicity, we have  $E(y_t) = E(\lambda_t) < \infty$  and hence  $E(y_t^k) < \infty$  and  $E(\lambda_t^k) < \infty$  for all  $k \in (0, 1]$ . For  $k > 1$ , we denote  $\tilde{y}_t = y_t^k$  and analogous to [22], we consider a sequence of perturbation  $\{\tilde{y}_t^m\}$  given by

$$\begin{cases} \tilde{y}_t^m = y_t^k + c_m \epsilon_t, & y_t \sim \text{Pois}(\lambda_t), \\ \log(\lambda_t) = \sum_{i=1}^{r+1} \left[ \alpha_i + \sum_{j=1}^{p_{Y,2}} \psi_{y,i}^{(j)} \log((\tilde{y}_{t-j}^m)^{1/k} + 1) \right] I((\tilde{y}_{t-4}^m)^{1/k} \in (\theta_{i-1}, \theta_i]), \end{cases}$$

where  $\{\epsilon_t\}$  are i.i.d. uniform random variables taking values on  $(0, 1)$ ,  $\{c_m\}$  is a sequence of positive constants converging to 0. Thus, we may utilize Theorem 7.1 for proving geometric ergodicity of  $\tilde{Y}_t^m = (\tilde{y}_t^m, \dots, \tilde{y}_{t-3}^m)$ , and show the existence of  $k$ th moment of  $\tilde{y}_t^m$ . As  $y_t$  follows Poisson distribution with mean  $\lambda_t$ , we have  $E(y_t^k | \lambda_t) \leq c_k \lambda_t^k + b_k$  for all  $k > 1$  and some positive constants  $c_k$  and  $b_k$ . Together with (7.17),

$$E(y_t^k | \tilde{Y}_{t-1}^m) \leq c_k E(\lambda_t^k | \tilde{Y}_{t-1}^m) + b_k \leq c_k (K_\lambda |\tilde{Y}_{t-1}^m|_\infty^{\rho/k} + c_\lambda)^k + b_k.$$

As  $\rho \in (0, 1)$ , for any constant  $c_1$  and  $c_2$ , we have  $c_1 y^\rho + c_2 < \rho y$  for all sufficiently large  $y$ . Therefore,

$$\begin{aligned} c_k (K_\lambda |\tilde{Y}_{t-1}^m|_\infty^{\rho/k} + c_\lambda)^k &= (c_k^{1/k} K_\lambda |\tilde{Y}_{t-1}^m|_\infty^{\rho/k} + c_k^{1/k} c_\lambda)^k \leq (\rho |\tilde{Y}_{t-1}^m|_\infty^{1/k})^k \\ &= \rho^k |\tilde{Y}_{t-1}^m|_\infty, \end{aligned}$$

for all  $|\tilde{Y}_{t-1}^m|_\infty^{1/k} > y^*$  with some sufficiently large  $y^*$ . Meanwhile, there exists a sufficiently large constant  $H_k$  such that  $c_k (K_\lambda |\tilde{Y}_{t-1}^m|_\infty^{\rho/k} + c_\lambda)^k \leq H_k$  for all  $|\tilde{Y}_{t-1}^m|_\infty^{1/k} \leq y^*$ . Hence,

$$E(y_t^k | \tilde{Y}_{t-1}^m) \leq \rho^k |\tilde{Y}_{t-1}^m|_\infty + H_k + b_\lambda = \rho^k |\tilde{Y}_{t-1}^m|_\infty + H_k^*,$$

where  $H_k^* = H_k + b_\lambda$ . Thus, we have

$$E(|\tilde{y}_t^m| | \tilde{Y}_{t-1}^m) = E(y_t^k | \tilde{Y}_{t-1}^m) + c_m E|\epsilon_t| \leq \rho^k |\tilde{Y}_{t-1}^m|_\infty + H_k^* + c_m/2.$$

By the induction arguments as in (7.12) and (7.13), we have

$$E(|\tilde{Y}_{t+3}^m|_\infty | \tilde{Y}_{t-1}^m) \leq (1 + \rho^k + \dots + \rho^{3k})(c_m/2 + H_k^*) + \rho^k |\tilde{Y}_{t-1}^m|_\infty. \tag{7.18}$$

Analogous to (7.13), (7.18) implies (7.2) with  $V(\tilde{Y}_{t-1}^m) = |\tilde{Y}_{t-1}^m|_\infty + 1$ , and hence the geometric ergodicity of  $\{\tilde{y}_t^m\}$  is verified by Theorem 7.1, which further implies  $E(\tilde{y}_t^m) < \infty$ . Then, as  $c_m \rightarrow 0$ ,  $\tilde{y}_t^m \rightarrow \tilde{y}_t$  almost surely and hence

$E(\tilde{y}_t) = E(y_t^k) < \infty$  for all  $k > 1$ . In conclusion,  $E(y_t^k) < \infty$  and  $E(\lambda_t^k) < \infty$  for all  $k > 0$ .

With the existence of all moments of  $y_t$  and  $\lambda_t$ , we are ready to verify the remaining conditions. As  $\log[E(e^{uT(y_t)} | \lambda_t)] = \lambda_t(e^u - 1)$ , by choosing  $c_t = \lambda_t(e^u - 1)$ , we have  $\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n c_t^{3/2} < \infty$ . Next, choosing

$$\Gamma(Y_{t-1}; y_t) = y_t \sum_{j=1}^2 \log(y_{t-j} + 1) + e^{K \sum_{j=1}^2 \log(y_{t-j} + 1)} \sum_{j=1}^2 \log(y_{t-j} + 1),$$

with some  $K > 0$ , we have  $E[\Gamma^2(Y_{t-1}; y_t)] < \infty$  since  $E(y_t^{2K})$  is finite for all  $K > 0$ . Therefore, all of the conditions in Assumption 4b) are satisfied.

**Proof of Theorems 5.1, 5.2, 5.3 and 5.4 (asymptotic theory of inferences)**

From (4.3), the minimum description length is a sum of the negative log-likelihood and the penalties on model complexity. We define

$$\begin{aligned} \text{Pen}(\Theta, p, d) &= \log_2(r) + \log_2(d) + \frac{1}{2} \sum_{i=1}^r \log_2(n_i) + \sum_{i=1}^{r+1} \log_2(p'_i + 4) \\ &\quad + \sum_{i=1}^{r+1} \frac{p'_i + 4}{2} \log_2(n_i), \end{aligned} \tag{7.19}$$

as the penalty on the model complexity. Similar to [66], we show the following propositions:

**Proposition 2.** *Let  $\theta_l < \theta_u$  be constants satisfying  $(\theta_l, \theta_u) \subseteq B_n = (\theta_{i-1}^0 - k_n^1, \theta_i^0 + k_n^2)$ , where  $\{k_n^1\}$  and  $\{k_n^2\}$  are two positive sequences that converging to 0. Define*

$$L_{n,i}(\Psi_i, \theta_l, \theta_u, d) = \frac{1}{n} \sum_{t=1}^n l(\Psi_i; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_l, \theta_u)).$$

*In addition, let  $L_i = \lim_{n \rightarrow \infty} E(L_{n,i})$ , and  $L_i^{(k)}, L_{n,i}^{(k)}$  be the  $k^{\text{th}}$  derivative with respect to  $\Psi_i$  for  $L_i, L_{n,i}$ , respectively. Here  $L_i^{(0)} = L_i, L_{n,i}^{(0)} = L_{n,i}$ . Then, under Assumptions 1-3,*

$$\sup_{\theta_l, \theta_u \in B_n} \sup_{\Psi_i} |L_i^{(k)}(\Psi_i, \theta_l, \theta_u, d) - L_{n,i}^{(k)}(\Psi_i, \theta_l, \theta_u, d)| \rightarrow 0 \quad \text{a.s.},$$

for  $k = 0, 1, 2, i = 1, \dots, r + 1$ .

*Proof of Proposition 2.* We show the case  $k = 0$  as an illustrative example. The proofs for  $k = 1$  or  $2$  are similar. First, by the compactness of the parameter space  $S(\Psi_i^*)$  of  $\Psi_i$  and the ergodic theorem, for any pair of  $(\theta_l, \theta_u) \subseteq [\theta_{i-1}^0, \theta_i^0]$ ,

$$\sup_{\Psi_i \in S(\Psi_i^*)} |L_i(\Psi_i, \theta_l, \theta_u, d) - L_{n,i}(\Psi_i, \theta_l, \theta_u, d)| \rightarrow 0 \quad \text{a.s.} \tag{7.20}$$

For discrete  $z_t$ , (7.20) holds for any combination of  $(\theta_l, \theta_u)$  from observed  $z_t$  values. Thus, Proposition 2 holds.

For continuous  $z_t$ , (7.20) holds particularly for all subintervals of rational endpoints. Hence for  $\theta \in (\theta_{i-1}^0, \theta_i^0)$  and any  $\epsilon > 0$ , there exists a rational number  $w < \theta$  such that

$$\begin{aligned} & \sup_{\Psi_i \in S(\Psi_i^*)} |L_i(\Psi_i, \theta_{i-1}^0, \theta, d) - L_{n,i}(\Psi_i, \theta_{i-1}^0, w, d)| \\ & \leq \sup_{\Psi_i \in S(\Psi_i^*)} \left| \frac{1}{n} \sum_{t=1}^n l(\Psi_i; y_t, Y_{t-1}, X_t) I(w < z_{t-d} \leq \theta) \right| + \epsilon/3 \\ & \leq \sup_{\Psi_i \in S(\Psi_i^*)} |E[l(\Psi_i; y_t, Y_{t-1}, X_t)] Q(w, \theta) + \epsilon/3| + \epsilon/3 \\ & < \epsilon, \end{aligned}$$

where  $Q(w, \theta) = E(I(w < z_{t-d} \leq \theta))$ . Selecting  $w$  close to  $\theta$  ensures a sufficiently small  $Q(w, \theta)$ . As the pair  $(\theta_l, \theta_u)$  is closely approximated by some subset with rational number endpoints, (7.20) holds uniformly on all  $(\theta_l, \theta_u) \subseteq [\theta_{i-1}^0, \theta_i^0]$ . Furthermore, if  $\theta_l$  or  $\theta_u$  is outside  $[\theta_{i-1}^0, \theta_i^0]$ , the almost sure convergence still holds if they are within a shrinkage neighborhood.  $\square$

**Proposition 3.** For any  $(\theta_l, \theta_u) \subseteq B_n = (\theta_{i-1}^0 - k_n^1, \theta_i^0 + k_n^2)$ , we define  $\Psi_i^* = \arg \max_{\Psi_i} L_i(\Psi_i, \theta_l, \theta_u, d)$  and  $\hat{\Psi}_i = \arg \max_{\Psi_i} L_{n,i}(\Psi_i, \theta_l, \theta_u, d)$ . Then, under Assumptions 1-3,  $\hat{\Psi}_i \rightarrow \Psi_i^*$  almost surely.

*Proof of Proposition 3.* From the definition of  $\hat{\Psi}_i$  and  $\Psi_i^*$ , we have

$$L_{n,i}(\hat{\Psi}_i, \theta_l, \theta_u, d) \geq L_{n,i}(\Psi_i^*, \theta_l, \theta_u, d), L_i(\hat{\Psi}_i, \theta_l, \theta_u, d) \leq L_i(\Psi_i^*, \theta_l, \theta_u, d).$$

By decomposing  $L_i(\Psi_i^*, \theta_l, \theta_u, d) - L_i(\hat{\Psi}_i, \theta_l, \theta_u, d)$ , we have almost surely that

$$\begin{aligned} 0 & \leq L_i(\Psi_i^*, \theta_l, \theta_u, d) - L_i(\hat{\Psi}_i, \theta_l, \theta_u, d) \\ & \leq |L_i(\Psi_i^*, \theta_l, \theta_u, d) - L_{n,i}(\Psi_i^*, \theta_l, \theta_u, d)| + L_{n,i}(\Psi_i^*, \theta_l, \theta_u, d) \\ & \quad - L_{n,i}(\hat{\Psi}_i, \theta_l, \theta_u, d) + |L_{n,i}(\hat{\Psi}_i, \theta_l, \theta_u, d) - L_i(\hat{\Psi}_i, \theta_l, \theta_u, d)| \quad (7.21) \\ & \leq 0, \end{aligned}$$

where the first and last terms in (7.21) converge to 0 almost surely by Proposition 2. Thus, by the uniqueness of  $\Psi_i^*$  under Assumption 2,  $\hat{\Psi}_i \rightarrow \Psi_i^*$  almost surely.  $\square$

For any  $(\theta_l, \theta_u) \subseteq (\theta_{i-1}^0, \theta_i^0)$ , by the theory of Kullback–Leibler distance,  $L_i(\Psi_i^0, \theta_l, \theta_u, d^0) \geq L_i(\Psi_i, \theta_l, \theta_u, d)$ , where the equality sign holds if and only if  $d = d^0$  and  $\Psi_i = \Psi_i^0$ . This observation is the main idea for proving Theorem 5.1.

*Proof of Theorem 5.1.* We prove by contradiction. Let  $A$  be the probability one set under which Propositions 2 and 3 hold. For each  $\omega \in A$ , suppose that there exists a subsequence  $\{n_m\}$  such that  $\hat{r}_{n_m} \rightarrow r^*$ ,  $\hat{d}_{n_m} \rightarrow d^*$ ,  $\hat{\Theta}_{n_m} \rightarrow \Theta^*$ ,  $\hat{\Psi}_{n_m} \rightarrow$



$\Psi^*$  almost surely with  $\{\Psi^*, \Theta^*, d^*\} \in \Omega_{\mathcal{M}} \times \{1, \dots, D\}$  under Assumption 1. For simplicity, we omit the subscript “ $n_m$ ” and replace  $\{n_m\}$  by  $\{n\}$  when necessary. Suppose that  $r^* < r^0$ , then there must be some  $\theta_{j-1}^*$  and  $\theta_j^*$  such that for some positive integer  $k$ ,

$$\theta_{i-1}^0 < \theta_{j-1}^* \leq \theta_i^0 < \theta_{i+1}^0 < \dots < \theta_{i+k}^0 < \theta_j^* \leq \theta_{i+k+1}^0.$$

In other words,  $k + 2$  true regimes are pooled into a “working regime”  $j$ . We relabel the true regimes as sub-regime 1, ...,  $k + 2$  with thresholds of sub-regime  $l$  denoted as  $(\theta_{(l-1)}, \theta_{(l)})$ . Hence, the number of observations in the working regime  $j$  is  $n_j$ . The log-likelihood of the working regime  $j$  is

$$n_j L_{n_j}^*(\hat{\Psi}_j, \hat{\theta}_{j-1}, \hat{\theta}_j, \hat{d}) = \sum_{l=1}^{k+2} n_j L_{n_j, l}^*(\hat{\Psi}_j, \hat{\theta}_{(l-1)}, \hat{\theta}_{(l)}, \hat{d}),$$

where

$$L_{n_j, l}^*(\Psi_j, \theta_{(l-1)}, \theta_{(l)}, d) = \frac{1}{n_j} \sum_{t=1}^{n_j} l(\Psi_j; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_{(l-1)}, \theta_{(l)}]),$$

$$L_{n_j}^*(\Psi_j, \theta_{j-1}^*, \theta_j^*, d) = \frac{1}{n_j} \sum_{t=1}^{n_j} l(\Psi_j; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_{j-1}^*, \theta_j^*]),$$

and  $\hat{\Psi}_j = \arg \max_{\Psi_j} L_{n_j}^*(\Psi_j, \theta_{j-1}^*, \theta_j^*, \hat{d})$ . Similarly, we define

$$L_{j, l}^*(\Psi_j, \theta_{(l-1)}, \theta_{(l)}, d) = \lim_{n_j \rightarrow \infty} E[L_{n_j, l}^*(\Psi_j, \theta_{(l-1)}, \theta_{(l)}, d)],$$

and

$$L_j^*(\Psi_j, \theta_{j-1}^*, \theta_j^*, d) = \lim_{n_j \rightarrow \infty} E[L_{n_j}^*(\Psi_j, \theta_{j-1}^*, \theta_j^*, d)].$$

As  $n \rightarrow \infty$ , we have  $n_j \rightarrow \infty$  and  $n_{j, l} \rightarrow \infty$ . Denote the true parameter in sub-regime  $l$  as  $\Psi_{(l)}^0$ . Consider a sub-regime  $m$ , from Propositions 2, 3 and theory of Kullback–Leibler distance,

$$\begin{aligned} \lim_{n \rightarrow \infty} L_{n_j, m}^*(\hat{\Psi}_j, \theta_{(m-1)}, \theta_{(m)}, \hat{d}) &= L_{j, m}^*(\Psi_j^*, \theta_{(m-1)}, \theta_{(m)}, d^*) \\ &\leq L_{j, m}^*(\Psi_{(m)}^0, \theta_{(m-1)}, \theta_{(m)}, d^0) \quad \text{a.s.} \end{aligned}$$

However, the equality cannot hold for all  $m$  under Assumption 2. Hence, we have

$$\lim_{n_j \rightarrow \infty} L_{n_j}^*(\hat{\Psi}_j, \hat{\theta}_{j-1}, \hat{\theta}_j, \hat{d}) < \sum_{l=1}^{k+2} L_{j, l}^*(\Psi_{(l)}^0, \theta_{(l-1)}, \theta_{(l)}, d^0) \quad \text{a.s.}, \quad (7.22)$$

since at least one part in the summation of (7.22) is not maximized. Furthermore, for one of such sub-regime  $m$ , equipped with the ergodic theorem, there exists some  $c_m > 0$  such that

$$\begin{aligned} L_{j,m}^*(\Psi_j^*, \theta_{(m-1)}, \theta_{(m)}, d^*) &= \mathbb{E}[l(\Psi_j^*; y_t, Y_{t-1}, X_t)I(z_{t-d^*} \in (\theta_{(m-1)}, \theta_{(m)}))] \\ &= \mathbb{E}[l(\Psi_{(m)}^0; y_t, Y_{t-1}, X_t)I(z_{t-d_0} \in (\theta_{(m-1)}, \theta_{(m)}))] \\ &\quad - c_m \\ &= L_{j,m}^*(\Psi_{(m)}^0, \theta_{(m-1)}, \theta_{(m)}, d^0) - c_m. \end{aligned}$$

By the ergodic theorem again, we have

$$L(\Psi^0, \Theta^0, d^0) - L(\Psi^*, \Theta^*, d^*) = O(n),$$

and is positive almost surely. On the other hand, from (7.19),  $\text{Pen}(\Theta^0, p^0, d^0) - \text{Pen}(\Theta^*, p^*, d^*)$  is of order  $O(\log(n))$ . Hence, the decrease in log-likelihood is more rapid. Therefore,  $r^* < r^0$  fails to optimize  $\text{MDL}(\mathcal{M})$ . In general, if one of the working regimes is not nested in a true regime, for example,  $\theta_{i-1}^0 < \theta_{j-1}^* < \theta_i^0 < \theta_j^* < \theta_{i+1}^0$  for some  $i, j$ , then the  $\text{MDL}(\mathcal{M})$  cannot be smaller than that with sub-regimes  $(\theta_{i-1}^0, \theta_{j-1}^*]$ ,  $(\theta_{j-1}^*, \theta_i^0]$  and  $(\theta_i^0, \theta_j^*]$ . As a result, all working regimes have to be nested in some true regimes, i.e.,  $\theta_{i-1}^0 < \theta_{j-1}^* < \theta_j^* < \theta_i^0$  for all  $j$  with some  $i$ , which implies  $r^* \geq r^0$ .

Next, assume that  $d^* \neq d^0$  with  $r^* \geq r^0$ . Consider regime  $i$ , from Proposition 2,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \sum_{i=1}^{r+1} l(\Psi_i; y_t, Y_{t-1}, X_t)I(z_{t-d} \in (\theta_{i-1}, \theta_i)) \\ = L_i(\Psi_i, \theta_{i-1}, \theta_i, d) < \infty \quad \text{a.s.} \end{aligned}$$

As the estimated regimes are nested in true regimes,  $(\theta_{j-1}^*, \theta_j^*) \subseteq [\theta_{i-1}^0, \theta_i^0]$ . By the property of Kullback-Leibler distance,

$$L_i(\Psi_i^0, \theta_{j-1}^*, \theta_j^*, d^0) - L_i(\Psi_j^*, \theta_{j-1}^*, \theta_j^*, d^*) \geq 0,$$

for all  $j$  and  $i$ , where the equality holds if and only if  $d^* = d^0$ . Hence it implies

$$L(\Theta^0, \Psi^0, d^0) > L(\Theta^*, \Psi^*, d^*), \tag{7.23}$$

where the difference between the two terms is of order  $O(n)$ . Thus, for the optimal model,  $\hat{d} \rightarrow d^0$  almost surely.

Moreover, if  $\hat{r} > r^0$ , at least two of the classified regimes are sub-regimes of a true regime. By Taylor's expansion, for any  $(\theta_l, \theta_u] \subseteq (\theta_{i-1}^0, \theta_i^0]$ ,

$$\begin{aligned} L_{n,i}(\hat{\Psi}_i, \theta_l, \theta_u, d^0) - L_{n,i}(\Psi_i^0, \theta_l, \theta_u, d^0) \\ = (\hat{\Psi}_i - \Psi_i^0)L_{n,i}^{(1)}(\Psi_i^0, \theta_l, \theta_u, d^0) + (\hat{\Psi}_i - \Psi_i^0)^2L_{n,i}^{(2)}(\Psi_i^0, \theta_l, \theta_u, d^0) + o(1) \\ = o(1) + (\hat{\Psi}_i - \Psi_i^0)^2L_{n,i}^{(2)}(\Psi_i^0, \theta_l, \theta_u, d^0) + o(1) \quad \text{a.s.} \end{aligned} \tag{7.24}$$

As

$$L_{n,i}^{(1)}(\hat{\Psi}_i, \theta_l, \theta_u, d^0) - L_{n,i}^{(1)}(\Psi_i^0, \theta_l, \theta_u, d^0) = (\hat{\Psi}_i - \Psi_i^0)L_{n,i}^{(2)}(\Psi_i^0, \theta_l, \theta_u, d^0) + o(1), \tag{7.25}$$

where the left-hand-side converges to 0 almost surely, by Assumption 2 and Kolmogorov’s law of iterated logarithm, we have

$$\begin{aligned} & \limsup_{n \rightarrow \infty} \frac{n^{1/2}}{(\log \log(n))^{1/2}} L_{n,i}^{(2)}(\Psi_i^0, \theta_l, \theta_u, d^0) \\ & \rightarrow 2^{1/2} [\text{var}(L_i^{(2)}(\Psi_i^0, \theta_l, \theta_u, d^0))]^{1/2} \quad \text{a.s.} \end{aligned} \quad (7.26)$$

Together with (7.25), we have  $|\hat{\Psi} - \Psi_i^0| = O((n^{-1} \log \log(n))^{1/2})$  almost surely. Thus, by (7.24),  $L(\Theta^0, \Psi^0, d^0) - L(\Theta^*, \Psi^*, d^*) = O(\log \log(n))$ . As  $r^* > r^0$  results in  $\text{Pen}(\Theta^*, p^*, d^0) > \text{Pen}(\Theta^0, p^0, d^0)$  where  $\text{Pen}(\Theta^*, p^*, d^0) - \text{Pen}(\Theta^0, p^0, d^0) = O(\log(n))$ , it implies that the MDL criterion is not minimized almost surely. Hence,  $\hat{r} \rightarrow r^0$  almost surely.

Under the consistency conditions, we can assume  $\hat{d} = d^0$  and  $\hat{r} = r^0$  holds for sufficiently large  $n$ . As all of the estimated regimes are nested in true regimes,  $\hat{\Theta} \rightarrow \Theta^0$  almost surely. From Proposition 2, in regime  $i$ , we have almost surely that  $\hat{\Psi}_i \rightarrow \Psi_i^*$ ,  $\hat{p}_i \rightarrow p_i^*$  for some  $\Psi_i^*$  and  $p_i^*$ . Suppose that  $\Psi_i^* \neq \Psi_i^0$  in one of the regime  $i$ . Then, as Kullback-Leibler distance suggests,  $L_i(\Psi_i^*, \theta_{i-1}^0, \theta_i^0, d^0) \leq L_i(\Psi_i^0, \theta_{i-1}^0, \theta_i^0, d^0)$ . Analogous to the argument in proving  $\hat{d} \rightarrow d^0$ , MDL( $\mathcal{M}$ ) is not asymptotically minimized. On the other hand, if  $\hat{\Psi} \rightarrow \Psi^0$  almost surely, we have  $\hat{p} \rightarrow p^*$  simultaneously for some  $p^*$ . As  $\text{Pen}(\Theta^0, p^*, d^0) \geq \text{Pen}(\Theta^0, p^0, d^0)$  where equality holds if and only if  $p^* = p^0$ , by the strong consistency of order estimation of MDL (see [47] and [48]) with respect to true regimes,  $\hat{p} \rightarrow p^0$  almost surely.  $\square$

*Proof of Theorem 5.2.*

1) We prove by showing the  $O_p(n^{-1})$  convergence speed of  $\hat{\theta}_i$  for all  $i$ . Under Theorem 5.1,  $|\hat{\theta}_i - \theta_i^0| = o(1)$  almost surely; hence it suffices to verify  $\text{pr}(|\hat{\theta}_i - \theta_i^0| > c/n) \rightarrow 0$  for some  $c > 0$ . Below we show  $\text{pr}(|\hat{\theta}_i - \theta_i^0| I(\theta_i^0 < \hat{\theta}_i) > c/n) \rightarrow 0$ , and the same argument give  $\text{pr}(|\hat{\theta}_i - \theta_i^0| I(\theta_i^0 \geq \hat{\theta}_i) > c/n) \rightarrow 0$ , and thus the result follows.

From the strong consistency of  $\hat{\theta}_i$ , we can restrict  $|\hat{\theta}_i - \theta_i^0| < \delta$  and  $|\hat{\Psi}_i - \Psi_i^0| < \delta$  for some positive  $\delta$  and sufficiently large  $n$ . Denote  $Q(a) = \text{E}\{I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a])\}$  for some  $a > 0$ , where  $I(\cdot)$  is the indicator function. Then for any  $\xi, \delta > 0$ , there is some  $c > 0$  such that

$$\text{pr} \left( \sup_{c/n < a \leq \delta} \left| \sum_{t=1}^n \frac{I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a])}{nQ(a)} - 1 \right| < \xi \right) > 1 - \epsilon. \quad (7.27)$$

Furthermore, from Assumption 4a), there exists a measurable function  $\Gamma_t = \Gamma(Y_{t-1}, X_t, y_t)$  such that  $\text{E}(\Gamma_t^2 | z_{t-d^0}) \leq M$  for all  $z_{t-d^0} \in [\theta_i^0 - \delta, \theta_i^0 + \delta]$  and some constant  $M > 0$ . Combined with Assumption 5, the joint process  $\{\Gamma_t I(z_{t-d^0} \in [\theta_i^0 - \delta, \theta_i^0 + \delta]), z_{t-d^0} I(z_{t-d^0} \in [\theta_i^0 - \delta, \theta_i^0 + \delta])\}$  is  $\rho$ -mixing with summable mixing coefficients. Therefore, we have

$$\begin{aligned} \text{pr} \left( \sup_{c/n < a \leq \delta} \left| \sum_{t=1}^n \frac{\Gamma_t I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a)) - \mathbb{E}\{\Gamma_t I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a))\}}{nQ(a)} \right| < \xi \right) \\ > 1 - \epsilon. \end{aligned} \quad (7.28)$$

The derivations of (7.27) and (7.28) follow similar arguments in Proposition 1 of [10] and Theorem 2 of [54]. Let

$$\begin{aligned} \mathcal{H}(\theta_i^0, a) &= \frac{1}{nQ(a)} \sum_{t=1}^n \left[ l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_{i+1}; y_t, Y_{t-1}, X_t) \right] \\ &\quad \times I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a)) \\ &= \frac{1}{nQ(a)} \sum_{t=1}^n \left[ \left\{ l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t) - l(\Psi_i^0; y_t, Y_{t-1}, X_t) \right\} \right. \\ &\quad + \left\{ l(\Psi_i^0; y_t, Y_{t-1}, X_t) - l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) \right\} \\ &\quad + \left. \left\{ l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_{i+1}; y_t, Y_{t-1}, X_t) \right\} \right] \\ &\quad \times I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a)). \end{aligned}$$

From Assumptions 4a) and 5, (7.27) and (7.28), we have for every  $\xi, \delta > 0$ , there exists a  $c > 0$  satisfying  $c/n < a < \delta$  such that

$$\begin{aligned} \mathcal{H}(\theta_i^0, a) &< \frac{1}{nQ(a)} \sum_{t=1}^n \left[ l(\Psi_i^0; y_t, Y_{t-1}, X_t) - l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) \right] \\ &\quad \times I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a)) \\ &\quad + \left( |\hat{\Psi}_i - \Psi_i^0| + |\Psi_{i+1}^0 - \hat{\Psi}_{i+1}| \right) (\xi + M) + \xi, \end{aligned} \quad (7.29)$$

with probability greater than  $1 - \epsilon$  for some small positive  $\epsilon$ . Using Assumptions 2, 4 and Lemma 5.35 in [62], for  $0 < a < \delta$ , there exists a constant  $\chi < 0$  such that

$$\mathbb{E} \left[ l(\Psi_i^0; y_t, Y_{t-1}, X_t) - l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) \mid z_{t-d^0} \in (\theta_i^0, \theta_i^0 + a) \right] \leq \chi < 0.$$

Following Theorem 2 of [54], the first term on the right hand side of (7.29) is less than  $\chi(1 - \xi)$  with some  $\xi > 0$ . Hence, by appropriate choices of  $\delta$  and  $\xi$  such that  $2\delta(M + \xi) + \xi + \chi(1 - \xi) < 0$ , we verify that  $\mathcal{H}(\theta_i^0, a) < 0$  in probability, with  $\delta, \xi$  and  $\epsilon$  go to 0 as  $n \rightarrow \infty$ .

Now suppose that  $|\theta_i^0 - \hat{\theta}_i| > c/n$  holds for some  $c > 0$  with positive probability. Denote  $\tilde{\Theta} = (\tilde{\theta}_1, \dots, \tilde{\theta})$  with  $\tilde{\theta}_j = \hat{\theta}_j$  if  $j \neq i$  and  $\tilde{\theta}_i = \theta_i^0$ , and  $\tilde{\Psi}$  is the maximum likelihood estimator under  $\tilde{\Theta}$ . As  $\hat{\Theta}$  is the minimizer of MDL, we have

$$\begin{aligned} &\text{pr} \left( |\theta_i^0 - \hat{\theta}_i| I(\theta_i^0 < \hat{\theta}_i) > c/n \right) \\ &= \text{pr} \left( |\theta_i^0 - \hat{\theta}_i| I(\theta_i^0 < \hat{\theta}_i) > c/n, \text{MDL}(\hat{\mathcal{M}}(\hat{\Psi}, \hat{\Theta}, d^0)) < \text{MDL}(\tilde{\mathcal{M}}(\tilde{\Psi}, \tilde{\Theta}, d^0)) \right) \\ &\leq \text{pr} \left( \text{MDL}(\hat{\mathcal{M}}(\hat{\Psi}, \hat{\Theta}, d^0)) < \text{MDL}(\tilde{\mathcal{M}}(\tilde{\Psi}, \tilde{\Theta}, d^0)) \right), \end{aligned}$$

where  $\hat{\mathcal{M}}(\hat{\Psi}, \hat{\Theta}, d^0)$ ,  $\tilde{\mathcal{M}}(\tilde{\Psi}, \tilde{\Theta}, d^0)$  are models with parameters  $\{\hat{\Psi}, \hat{\Theta}, d^0\}$  and  $\{\tilde{\Psi}, \tilde{\Theta}, d^0\}$ , respectively. As  $|\theta_i^0 - \hat{\theta}_i| > c/n$ , there exists some  $a^* \in (c/n, \delta]$  with a sufficiently small  $\delta$  such that

$$L(\tilde{\Psi}, \tilde{\Theta}, d^0) - L(\hat{\Psi}, \hat{\Theta}, d^0) = -nQ(a^*)\mathcal{H}(\theta_i^0, a^*) = O_p(1). \tag{7.30}$$

Furthermore,  $\mathcal{H}(\theta_i^0, a)$  is negative with probability going to 1, the term in (7.30) converges in probability to a positive term. Meanwhile, as  $\hat{\theta}_i \rightarrow \theta_i^0$  and  $\hat{p}_i \rightarrow p_i^0$  almost surely, the difference of penalties  $\text{Pen}(\tilde{\Theta}, \tilde{p}, d^0) - \text{Pen}(\hat{\Theta}, \hat{p}, d^0) = o(1)$ . Therefore, as  $n \rightarrow \infty$ ,

$$\text{pr}(\text{MDL}(\hat{\mathcal{M}}(\hat{\Psi}, \hat{\Theta}, d^0)) > \text{MDL}(\tilde{\mathcal{M}}(\tilde{\Psi}, \tilde{\Theta}, d^0))) \rightarrow 1,$$

and thus  $\text{pr}(|\theta_i^0 - \hat{\theta}_i| I(\theta_i^0 < \hat{\theta}_i) > c/n) \rightarrow 0$ . Therefore, the convergence rate of  $\hat{\theta}_i$  is of order  $O_p(n^{-1})$  for all  $i = 1, \dots, r^0$ , and the proof is complete.

2) We first prove the convergence property of  $\text{pr}(\hat{\theta}_i \neq \theta_i^0)$ . For simplicity we assume  $\hat{\theta}_i$  are integers and  $\nu_i = \nu$  in all regimes. Without loss of generality, we show the convergence of  $\text{pr}(\hat{\theta}_i > \theta_i^0)$ .

Denote  $\hat{\Psi}_i = \arg \max_{\Psi_i} L_{n,i}(\Psi_i, \theta_{i-1}^0, \theta_i^0, d^0)$  and  $\hat{\Psi}_i^* = \arg \max_{\Psi_i} L_{n,i}(\Psi_i, \theta_{i-1}^0, \tilde{\theta}_i, d^0)$ , where  $\hat{\Psi}_i^* \rightarrow \Psi_i^*$  almost surely for  $\Psi_i^* = \arg \max_{\Psi_i} L_{n,i}(\Psi_i, \theta_{i-1}^0, \tilde{\theta}_i, d^0)$  which is neither  $\Psi_i^0$  nor  $\Psi_{i+1}^0$ . Under (5.2), we have  $\lambda_t^* = \lambda_t(\Psi_i^*; Y_{t-1}, X_t, z_{t-d})$  and  $\lambda_t^0 = \lambda_t(\Psi_i^0; Y_{t-1}, X_t, z_{t-d})$ , and  $\lambda_t^* \neq \lambda_t^0$  almost surely.

For discrete  $z_t$  and any  $\tilde{\theta}_i > \theta_i^0$ , the number of observations with  $z_{t-d} \in (\theta_i^0, \tilde{\theta}_i]$  with order  $O(n)$  almost surely. If  $\hat{\theta}_i \rightarrow \tilde{\theta}_i \neq \theta_i^0$ , then by the ergodic theorem,

$$n \left[ L_{n,i}(\Psi_i^0, \theta_{i-1}^0, \theta_i^0, d^0) - L_{n,i}(\Psi_i^*, \theta_{i-1}^0, \tilde{\theta}_i, d^0) \right] = O(n) \quad \text{a.s.},$$

which overweights the  $O(\log(n))$  rate of difference in penalties under  $\theta_i^0$  and  $\tilde{\theta}_i$  by (7.19). Hence we claim that asymptotically,

$$\begin{aligned} & \text{pr}(\hat{\theta}_i > \theta_i^0) \\ = & \text{pr} \left( \sum_{t=1}^n [l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t)] I(z_{t-d} \in (\theta_{i-1}^0, \hat{\theta}_i]) > 0 \right). \end{aligned}$$

For any  $\tilde{\theta}_i > \theta_i^0$ , we divide all observations with  $z_{t-d} \in (\theta_{i-1}^0, \tilde{\theta}_i]$  into two partitions: the first partition constitutes all  $y_t$  with  $z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]$ ; and the second partition constitutes of all  $y_t$  with  $z_{t-d} \in (\theta_i^0, \tilde{\theta}_i]$ . For the first partition, from Assumption 4b),

$$\begin{aligned} & l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t) \\ = & [l(\Psi_i^*; y_t, Y_{t-1}, X_t) - l(\Psi_i^0; y_t, Y_{t-1}, X_t)](1 + o(1)) \\ = & \frac{1}{\nu a_t} [\{T(y_t)\gamma(\lambda_t^*) - b(\lambda_t^*)\} - \{T(y_t)\gamma(\lambda_t^0) - b(\lambda_t^0)\}](1 + o(1)). \end{aligned} \tag{7.31}$$

Here  $\lambda_t^*$  and  $\lambda_t^0$  are deterministic given  $\{Y_{t-1}, X_t, z_{t-d}\}$ , and  $\{T(y_t)\}_{t=1,2,\dots}$  in (7.31) are mutually independent given  $\lambda_t^0$ . To simplify derivation, we relabel all observations in the first group from  $s = 1, \dots, n_i^0$ , and denote  $\Lambda_{n_i^0}$  as a realization of  $\{\lambda_s^*, \lambda_s^0\}_{s=1,\dots,n_i^0}$ , where  $\lambda_s^* \neq \lambda_s^0$ ,  $s = 1, \dots, n_i^0$  almost surely. Let  $f_\Lambda$  be the joint density of  $\{\lambda_s^*, \lambda_s^0\}_{s=1,\dots,n_i^0}$ . By (7.31), for sufficiently large  $n$  and  $n_i^0$ ,

$$\begin{aligned} & \text{pr} \left( \sum_{t=1}^n [l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t)] I(z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]) > 0 \right) \\ &= \int \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} [\{T(y_s)\gamma(\lambda_s^*) - b(\lambda_s^*)\} \right. \\ &\quad \left. - \{T(y_s)\gamma(\lambda_s^0) - b(\lambda_s^0)\}](1 + o(1)) > 0 \mid \Lambda_{n_i^0} \right) \\ &\quad \times f_\Lambda \left( \{\lambda_s^*, \lambda_s^0\}_{s=1,\dots,n_i^0} \right) d\lambda_1^* \lambda_1^0 \dots \lambda_{n_i^0}^* \lambda_{n_i^0}^0, \quad (7.32) \end{aligned}$$

with

$$\begin{aligned} & \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} [\{T(y_s)\gamma(\lambda_s^*) - b(\lambda_s^*)\} \right. \\ &\quad \left. - \{T(y_s)\gamma(\lambda_s^0) - b(\lambda_s^0)\}](1 + o(1)) > 0 \mid \Lambda_{n_i^0} \right) \\ &= \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} [T(y_s)\{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\} - E(T(y_s) \mid \lambda_s^0)\{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}] \right. \\ &\quad \left. > \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} [b(\lambda_s^*) - b(\lambda_s^0) - E(T(y_s) \mid \lambda_s^0)\{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}] \mid \Lambda_{n_i^0} \right) \\ &\quad \times (1 + o(1)) \\ &= \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} T^*(y_s)\{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\} > \sum_{s=1}^{n_i^0} \Delta_b(\lambda_s^*, \lambda_s^0) \mid \Lambda_{n_i^0} \right) (1 + o(1)), \quad (7.33) \end{aligned}$$

where given  $\Lambda_{n_i^0}$  in (7.33),  $T^*(y_s) = T(y_s) - E(T(y_s) \mid \lambda_s^0)$ ,  $s = 1, 2, \dots$  are mutually independent. Denote

$$\begin{aligned} \Delta_b(\lambda_s^*, \lambda_s^0) &= \frac{1}{\nu a_s} [b(\lambda_s^*) - b(\lambda_s^0) - E(T(y_s) \mid \lambda_s^0)\{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}] \\ &= E[l(\Psi_i^0; y_s, Y_{s-1}, X_s)] - E[l(\Psi_i^*; y_t, Y_{s-1}, X_s)]. \end{aligned}$$

By Assumption 4b),  $E[l(\Psi_i^0; y_s, Y_{s-1}, X_s) - l(\Psi_i^*; y_t, Y_{s-1}, X_s)]^2 < \infty$ , and hence

$$E|\Delta_b(\lambda_s^*, \lambda_s^0)| < \infty, \quad \text{and} \quad E\left[\frac{1}{\nu^2 a_s^2} \text{var}(T(y_s) \mid \lambda_s^0) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}^2\right] < \infty. \tag{7.34}$$

By the property of Kullback–Leibler distance,  $E[l(\Psi_i^0; y_s, Y_{s-1}, X_s)] > E[l(\Psi_i^*; y_t, Y_{s-1}, X_s)]$  if  $\lambda_s^* \neq \lambda_s^0$ . Thus  $\Delta_b(\lambda_s^*, \lambda_s^0) > 0$  for all  $y_s$  with  $z_{s-d} \in (\theta_{i-1}^0, \theta_i^0]$ . As  $\lambda_s^* \neq \lambda_s^0$  almost surely and  $\{Y_{t-1}, X_t, z_{t-d}\}$  is ergodic, by (7.34) and the ergodic theorem,

$$\lim_{n_i^0 \rightarrow \infty} \frac{1}{n_i^0} \sum_{s=1}^{n_i^0} \Delta_b(\lambda_s^*, \lambda_s^0) = E(\Delta_b(\lambda_s^*, \lambda_s^0)) > 0 \quad \text{a.s.} \tag{7.35}$$

Meanwhile, as

$$\begin{aligned} & \text{var} \left[ \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} T^*(y_s) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\} \mid \Lambda_{n_i^0} \right] \\ &= \sum_{s=1}^{n_i^0} \frac{1}{\nu^2 a_s^2} \text{var}(T(y_s) \mid \lambda_s^0) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}^2, \end{aligned}$$

analogously we have

$$\begin{aligned} & \lim_{n_i^0 \rightarrow \infty} \frac{1}{n_i^0} \text{var} \left[ \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} T^*(y_s) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\} \mid \Lambda_{n_i^0} \right] \\ &= \lim_{n_i^0 \rightarrow \infty} \frac{1}{n_i^0} \sum_{s=1}^{n_i^0} \frac{1}{\nu^2 a_s^2} \text{var}(T(y_s) \mid \lambda_s^0) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}^2 \\ &= E \left[ \frac{1}{\nu^2 a_s^2} \text{var}(T(y_s) \mid \lambda_s^0) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}^2 \right] > 0 \quad \text{a.s.} \tag{7.36} \end{aligned}$$

by (7.34) and the ergodic theorem. Denote

$$\mathcal{X}_{n_i^0} = \frac{\sum_{s=1}^{n_i^0} \Delta_b(\lambda_s^*, \lambda_s^0)}{\left[ \sum_{s=1}^{n_i^0} \text{var} \left( \frac{1}{\nu a_s} T^*(y_s) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\} \mid \Lambda_{n_i^0} \right) \right]^{1/2}}. \tag{7.37}$$

By (7.35) and (7.36),  $(n_i^0)^{-1/2} \mathcal{X}_{n_i^0}$  converges almost surely to some  $\tau = \tau(\tilde{\theta}_i) > 0$ . By Theorem 10 in Chapter VIII of [46] and Assumption 4b), we have

$$\begin{aligned} & \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} T^*(y_s) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\} > \sum_{s=1}^{n_i^0} \Delta_b(\lambda_s^*, \lambda_s^0) \mid \Lambda_{n_i^0} \right) \\ &= \{1 - \Phi(\mathcal{X}_{n_i^0})\} \exp \left[ \frac{\mathcal{X}_{n_i^0}^3}{(n_i^0)^{1/2}} h_{n_i^0} \left( \frac{\mathcal{X}_{n_i^0}}{(n_i^0)^{1/2}} \right) \right] (1 + l_1 \tau), \tag{7.38} \end{aligned}$$

where  $\Phi(\cdot)$  is the cumulative distribution function of a standard normal random variable,  $h_{n_i^0}(x) = \sum_{k=0}^{\infty} a_{kn_i^0} x^k$  is the generalized Cramér series defined on p.220 of [46], and  $l_1$  is some constant. By [46],  $h_{n_i^0}(\tau)$  converges as  $n_i^0$  goes to infinity for sufficiently small  $|\tau|$ . Hence, we may select a sufficiently small constant  $\tau^*$  such that for sufficiently large  $n_i^0$ ,  $h_{n_i^0}(\tau) \leq H$  for some constant  $H$  whenever  $|\tau| \leq \tau^*$ . Therefore, for sufficiently large  $n$ , if  $\tau \leq \tau^*$ , we have

$$1 - \Phi(\mathcal{X}_{n_i^0}) \leq c_1 n^{-\frac{1}{2}} e^{-\frac{1}{2}\tau^2 n}, \quad \text{and} \quad \exp \left[ \frac{\mathcal{X}_{n_i^0}^3}{(n_i^0)^{1/2}} h_{n_i^0} \left( \frac{\mathcal{X}_{n_i^0}}{(n_i^0)^{1/2}} \right) \right] \leq c_2 e^{\tau^3 n},$$

for some  $c_1, c_2 > 0$ . Thus, (7.38) is of order  $O(n^{-1/2} e^{-an})$  almost surely for some  $a = \tau^2/2 - \tau^3 > 0$ , by choosing a proper  $\tau^* < 1/2$ . Otherwise, if  $\tau > \tau^*$ , we have

$$\begin{aligned} & \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} T^*(y_s) \{ \gamma(\lambda_s^*) - \gamma(\lambda_s^0) \} > \sum_{s=1}^{n_i^0} \Delta_b(\lambda_s^*, \lambda_s^0) \mid \Lambda_{n_i^0} \right) \\ & \leq \text{pr} \left( \sum_{s=1}^{n_i^0} \frac{1}{\nu a_s} T^*(y_s) \{ \gamma(\lambda_s^*) - \gamma(\lambda_s^0) \} > \frac{\tau^*}{\tau} \sum_{s=1}^{n_i^0} \Delta_b(\lambda_s^*, \lambda_s^0) \mid \Lambda_{n_i^0} \right) \\ & = \left\{ 1 - \Phi \left( \frac{\tau^*}{\tau} \mathcal{X}_{n_i^0} \right) \right\} \exp \left[ \frac{(\tau^*)^3}{\tau^3} \frac{\mathcal{X}_{n_i^0}^3}{(n_i^0)^{1/2}} h_{n_i^0} \left( \frac{\tau^* \mathcal{X}_{n_i^0}}{\tau (n_i^0)^{1/2}} \right) \right] (1 + l_1 \tau^*) \\ & \leq c_1 c_2 (1 + l_1 \tau^*) n^{-\frac{1}{2}} e^{-\left(\frac{1}{2}(\tau^*)^2 - (\tau^*)^3\right)n} \quad \text{a.s.} \end{aligned} \tag{7.39}$$

Next, using similar arguments in (7.38), (7.39) is of order  $O(n^{-1/2} e^{-a'n})$  almost surely for some  $a' = (\tau^*)^2/2 - (\tau^*)^3 > 0$ . Thus, given  $\Lambda_{n_i^0}$ , (7.33) is of order  $O(n^{-1/2} e^{-a_1 n})$  almost surely for some  $a_1 > 0$ . Combining with (7.32), we have

$$\begin{aligned} & \text{pr} \left( \sum_{t=1}^n [l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t)] I(z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]) > 0 \right) \\ & = O(n^{-1/2} e^{-a_1 n}). \end{aligned}$$

Applying similar argument for the second partition, we have

$$\begin{aligned} & \text{pr} \left( \sum_{t=1}^n [l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_{i+1}; y_t, Y_{t-1}, X_t)] I(z_{t-d} \in (\theta_i^0, \tilde{\theta}_i]) > 0 \right) \\ & = O(n^{-1/2} e^{-a_2 n}), \end{aligned}$$

for some  $a_2 > 0$ . Hence,

$$\begin{aligned} & \text{pr}(\hat{\theta}_i = \tilde{\theta}_i) \\ & \leq \text{pr} \left( \sum_{t=1}^n [l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t)] I(z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]) > 0 \right) \end{aligned}$$



$$\begin{aligned}
 & + \text{pr} \left( \sum_{t=1}^n [l(\hat{\Psi}_i^*; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_{i+1}; y_t, Y_{t-1}, X_t)] I(z_{t-d} \in (\theta_{i-1}^0, \tilde{\theta}_i]) > 0 \right) \\
 & = O(n^{-1/2} e^{-a_1 n}) + O(n^{-1/2} e^{-a_2 n}) = O(n^{-1/2} e^{-\min(a_1, a_2) n}). \tag{7.40}
 \end{aligned}$$

Note that  $a_1$  and  $a_2$  in (7.40) are functions of  $\tilde{\theta}_i$ . As the choices of  $\tilde{\theta}_i$  are infinite in regime 1 and  $r + 1$ , we need to show that  $a_1$  and  $a_2$  are positive uniform in all  $\tilde{\theta}_i$ . Without loss of generality we prove this for regime  $r + 1$ : we verify both  $\inf_{\tilde{\theta}_i} a_1(\tilde{\theta}_i) > 0$  and  $\inf_{\tilde{\theta}_i} a_2(\tilde{\theta}_i) > 0$ . From (7.39), showing  $\inf_{\tilde{\theta}_i} a_1 > 0$  is equivalent to showing  $\inf_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \tau(\tilde{\theta}_i) > 0$ . In addition, from (7.37), it suffices to show

$$\inf_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \mathbb{E}(\Delta_b(\lambda_s^*, \lambda_s^0)) > 0,$$

and

$$\sup_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \mathbb{E}[\text{var}(T(y_t) \mid \lambda_s^0) (\gamma(\lambda_s^*) - \gamma(\lambda_s^0))^2 / (va_s)^2] < \infty.$$

First, by Assumption 2, we have

$$\begin{aligned}
 & \mathbb{E}\{l(\Psi_i^*; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_{i-1}^0, \tilde{\theta}_i])\} \\
 & < \mathbb{E}\{l(\Psi_i^0; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_{i-1}^0, \theta_i^0])\} \\
 & \quad + \mathbb{E}\{l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) I(z_{t-d} \in (\theta_i^0, \tilde{\theta}_i])\}.
 \end{aligned}$$

Therefore, for all  $z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]$ , there exists some  $\beta \in (0, 1)$  such that

$$\begin{aligned}
 \mathbb{E}(l(\Psi_i^*; y_t, Y_{t-1}, X_t)) & \leq (1 - \beta) \mathbb{E}\{l(\Psi_i^0; y_t, Y_{t-1}, X_t)\} \\
 & \quad + \beta \mathbb{E}\{l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t)\}. \tag{7.41}
 \end{aligned}$$

Given  $z_{t-d} \in (\theta_{i-1}^0, \theta_i^0]$ , by (7.41) and the property of Kullback-Leibler distance, we have

$$\begin{aligned}
 & \mathbb{E}\{l(\Psi_i^0; y_t, Y_{t-1}, X_t)\} - \mathbb{E}\{l(\Psi_i^*; y_t, Y_{t-1}, X_t)\} \\
 & \geq \mathbb{E}\{l(\Psi_i^0; y_t, Y_{t-1}, X_t)\} \\
 & \quad - [(1 - \beta) \mathbb{E}\{l(\Psi_i^0; y_t, Y_{t-1}, X_t)\} + \beta \mathbb{E}\{l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t)\}] \\
 & = \beta [\mathbb{E}\{l(\Psi_i^0; y_t, Y_{t-1}, X_t)\} - \mathbb{E}\{l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t)\}] > 0. \tag{7.42}
 \end{aligned}$$

By (7.42),  $\inf_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \mathbb{E}(\Delta_b(\lambda_s^*, \lambda_s^0)) > 0$ . On the other hand, by Assumption 4b), we have

$$\begin{aligned}
 & \sup_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \mathbb{E} \left[ \frac{1}{\nu^2 a_s^2} \text{var}(T(y_s) \mid \lambda_s^0) \{\gamma(\lambda_s^*) - \gamma(\lambda_s^0)\}^2 \right] \\
 & = \sup_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \mathbb{E} \left[ \text{var} \left\{ l(\Psi_i^*; y_t, Y_{t-1}, X_t) - l(\Psi_i^0; y_t, Y_{t-1}, X_t) \mid Y_{t-1}, X_t \right\} \right] \\
 & \leq \sup_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} \mathbb{E} \left[ \mathbb{E} \left\{ l(\Psi_i^*; y_t, Y_{t-1}, X_t) - l(\Psi_i^0; y_t, Y_{t-1}, X_t) \mid Y_{t-1}, X_t \right\}^2 \right]
 \end{aligned}$$

$$\begin{aligned} &\leq \mathbb{E} \left[ \mathbb{E} \left\{ \Gamma^2(Y_{t-1}, X_t, y_t \mid Y_{t-1}, X_t) \right\} \sup_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} |\Psi_i^* - \Psi_i^0|^2 \right] \\ &= \mathbb{E}[\Gamma^2(Y_{t-1}, X_t, y_t)] \sup_{\tilde{\theta}_i \in [\theta_i^0 + 1, \theta_{i+1}^0]} |\Psi_i^* - \Psi_i^0|^2, \end{aligned} \tag{7.43}$$

where (7.43) is bounded by Assumption 4b) and the compactness of  $\Omega_\Psi$  in Assumption 1. Therefore,  $\inf_{\tilde{\theta}_i \in [\theta_i^0, \theta_{i+1}^0]} \tau > 0$  and hence  $\inf_{\tilde{\theta}_i \in [\theta_i^0, \theta_{i+1}^0]} a_1(\tilde{\theta}_i) > 0$ . Analogously,  $\inf_{\tilde{\theta}_i \in [\theta_i^0, \theta_{i+1}^0]} a_2(\tilde{\theta}_i) > 0$ . Since  $\hat{\theta}_i$  takes value in  $\{z_{t-d}\}$ , the number of choices in  $\tilde{\theta}_i$  is  $O(n)$  increasing. By (7.40),

$$\text{pr}(\hat{\theta}_i > \theta_i^0) = \sum_{\tilde{\theta}_i > \theta_i^0} \text{pr}(\hat{\theta}_i = \tilde{\theta}_i) = O(n^{1/2}e^{-a^*n}),$$

for some  $a^* > 0$ . The derivation of the same convergence rate of  $\text{pr}(\hat{\theta}_i < \theta_i^0)$  is analogous, and hence  $\text{pr}(\hat{\theta}_i \neq \theta_i^0) = O(n^{1/2}e^{-a_2^*n})$  for some  $a_2^* > 0$ . If  $\nu_i \neq \nu_{i+1}$  for some  $i$ , the derivation is similar albeit more tedious. Finally, as  $\text{pr}(\hat{\Theta} \neq \Theta^0) \leq \sum_{i=1}^r \text{pr}(\hat{\theta}_i \neq \theta_i^0)$ , the result in 2) follows.  $\square$

*Proof of Theorem 5.3.* First we prove the weak convergence of  $n(\hat{\theta}_i - \theta_i^0)$  to  $M_i^-$ , where  $[M_i^-, M_i^+)$  is the unique random interval which minimizes the process

$$\begin{aligned} \ddot{\ell}_i(\kappa_i) &= \sum_{t=1}^n [l(\hat{\Psi}_{i+1}; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t)] \\ &\quad \times I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + \frac{\kappa_i}{n}]) I(\kappa_i \geq 0) \\ &\quad + \sum_{t=1}^n [l(\hat{\Psi}_i; y_t, Y_{t-1}, X_t) - l(\hat{\Psi}_{i+1}; y_t, Y_{t-1}, X_t)] \\ &\quad \times I(z_{t-d^0} \in (\theta_i^0 + \frac{\kappa_i}{n}, \theta_i^0]) I(\kappa_i < 0), \end{aligned} \tag{7.44}$$

where  $|\hat{\theta}_i - \theta_i^0| = |\kappa_i|/n$ . When  $\hat{\Psi}_i$  and  $\hat{\Psi}_{i+1}$  are replaced by  $\Psi_i^0$  and  $\Psi_{i+1}^0$  in (7.44), we define an analogous process  $\tilde{\ell}_i(\kappa_i) = \tilde{\ell}_{1,i}(\kappa_i)I(\kappa_i \geq 0) + \tilde{\ell}_{2,i}(\kappa_i)I(\kappa_i < 0)$  where

$$\begin{aligned} \tilde{\ell}_{1,i}(\kappa_i) &= \sum_{t=1}^n [l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) - l(\Psi_i^0; y_t, Y_{t-1}, X_t)] I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + \frac{\kappa_i}{n}]), \\ \tilde{\ell}_{2,i}(\kappa_i) &= \sum_{t=1}^n [l(\Psi_i^0; y_t, Y_{t-1}, X_t) - l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t)] I(z_{t-d^0} \in (\theta_i^0 + \frac{\kappa_i}{n}, \theta_i^0]), \end{aligned}$$

which correspond to the case  $\hat{\theta}_i \geq \theta_i^0$  and  $\hat{\theta}_i < \theta_i^0$ , respectively. Denote  $\hat{\Psi}^{(0)} = (\hat{\Psi}_1^{(0)}, \dots, \hat{\Psi}_{r+1}^{(0)})$  as the maximum likelihood estimates given  $\Theta^0$ . By Assumptions 1–5 and Lemmas 3 and 4 of [54], we have for all  $K > 0$ ,

$$\sup_{|\hat{\theta}_i - \theta_i^0| < K/n} |\hat{\Psi}_i - \hat{\Psi}_i^{(0)}| = o_p(n^{-1/2}), \tag{7.45}$$

and

$$\sup_{\kappa_i \leq K} |\ddot{\ell}(\kappa_i) - \tilde{\ell}(\kappa_i)| = o_p(1). \tag{7.46}$$

By (7.46), it suffices to study the convergence properties of  $\tilde{\ell}_i(\kappa_i)$ . We show that  $\tilde{\ell}_{1,i}(\kappa_i)$  converges to  $\tilde{\ell}_{1,i}^*(\kappa_i)$  in distribution, and the convergence of  $\tilde{\ell}_{2,i}(\kappa_i)$  follows similarly. We study the right-continuous version of  $\tilde{\ell}_{1,i}^*(\kappa_i)$  with assuming  $\kappa_i \geq 0$ .

First, we construct a piecewise constant interpolation for the difference in log-likelihood. Taking  $\epsilon = 1/n$ , we let  $\eta^\epsilon(v) = X_{\lfloor nv \rfloor}^\epsilon$  for  $0 \leq v \leq 1$ , where  $\lfloor \cdot \rfloor$  is the floor function. Note that  $X_0^\epsilon = 0, X_{t+1}^\epsilon = X_t^\epsilon + \zeta_{t+1}^\epsilon$ . In addition, define

$$\zeta_t^\epsilon = [l(\Psi_{i+1}^0; y_t, Y_{t-1}, X_t) - l(\Psi_i^0; y_t, Y_{t-1}, X_t)]I(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + \kappa_i \epsilon]).$$

By construction, we have  $\eta^\epsilon(1) = \tilde{\ell}_{1,i}(\kappa_i)$ , and  $\eta^\epsilon(v) = X_s^\epsilon$  when  $v \in [s\epsilon, s\epsilon + \epsilon)$ . Using similar arguments for the discontinuities of  $\tilde{\ell}_{1,i}(\kappa_i)$  on  $(u, u+h]$  in the proof of Lemma 3.2 in [25],  $\{\tilde{\ell}_{1,i}(\kappa_i)\}$  is tight with  $\kappa_i \geq 0$ . Combining with the uniform boundedness of  $\eta^\epsilon(v)$  which is verified by the truncation argument in [29], we can apply the operator convergence in [29], by which the weak convergence of  $\eta^\epsilon(v)$  can be deduced. Then taking  $v = 1$ , the weak convergence of  $\tilde{\ell}_{1,i}(\kappa_i)$  in the Skorohod metric is established.

Next, we show the weak convergence of  $\eta^\epsilon(v)$  to the compound Poisson process  $\mathcal{C}(v)$  with intensity  $\pi(\theta_i^0)\kappa$ . Let  $\mathcal{F}_v$  be a sequence of  $\sigma$ -algebras that are generated from  $\{\eta^\epsilon(u), u \leq v\}$ , and  $\mathcal{L}$  be a set of progressively measurable functions  $f$  with respect to  $\mathcal{F}_v$ , with  $\sup_v E|f(v)| < \infty$ . Then, let  $\mathcal{F}_v^\epsilon \subseteq \mathcal{F}_v$  denote the  $\sigma$ -algebra generated by  $\{\eta^\epsilon(u), u \leq v\}$ , and  $E_v^\epsilon$  denote the conditional expectation under  $\mathcal{F}_v^\epsilon$ . From the definition in [28], we denote  $p\text{-lim}_{\delta \rightarrow 0} f^\delta = f$  for  $f^\delta \in \mathcal{L}$  if and only if  $\lim_{\delta \rightarrow 0} E|f^\delta - f| = 0$ . In addition, define the  $p$ -infinitesimal operator  $\hat{A}^\epsilon$  by  $\hat{A}^\epsilon f^\epsilon = p\text{-lim}_{\delta \rightarrow 0} [E_v^\epsilon f(v + \delta) - f(v)]/\delta$ , which is  $p$ -right continuous and in  $\mathcal{L}$ . Denote  $\hat{\mathcal{L}}$  as the space of continuous bounded positive function  $f$  with  $\lim_{v \rightarrow \infty} f(v) \rightarrow 0$ , and  $\hat{\mathcal{L}}^{(2)} \subseteq \hat{\mathcal{L}}$  with compact supports and continuous bounded second derivatives. Then, define the operator  $A$  on  $\hat{\mathcal{L}}^{(2)}$  as

$$Af(v) = \pi_z(\theta_i^0)\kappa_i \int (f(y+v) - f(v))q(dy),$$

where  $q(dy)$  is the induced probability measure of  $l(y_t; Y_{t-1}, X_t, \Psi_{i+1}^0) - l(y_t; Y_{t-1}, X_t, \Psi_i^0)$  conditioning on  $z_{t-d^0} = (\theta_i^0)^+$ . Let  $f \in \hat{\mathcal{L}}^{(2)}$ , and for any  $\tau^\epsilon > 0$ , define

$$f^\epsilon(v) = \frac{1}{\tau^\epsilon} \int_0^{\tau^\epsilon} E_v^\epsilon [f(\eta^\epsilon(v+s))] ds.$$

By construction,  $f^\epsilon(v) \in \hat{\mathcal{L}}^{(2)}$  and  $p\text{-lim}_{\tau^\epsilon \rightarrow 0} f^\epsilon(v) = f(\eta^\epsilon(v))$ . Moreover, from [28],

$$\hat{A}^\epsilon f^\epsilon(v) = \frac{1}{\tau^\epsilon} [E_v^\epsilon \{f(\eta^\epsilon(v + \tau^\epsilon))\} - f(\eta^\epsilon(v))].$$

Next we show the operator convergence of  $\hat{A}^\epsilon f^\epsilon(v)$  to  $Af(\eta^\epsilon(v))$ . Express

$$\hat{A}^\epsilon f^\epsilon(v) = \frac{1}{\tau^\epsilon} \sum_{k=0}^{\lfloor n(v+\tau^\epsilon) \rfloor - \lfloor nv \rfloor - 1} \mathbb{E}_v^\epsilon \{ f(\eta^\epsilon(v+k\epsilon) + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(\eta^\epsilon(v+k\epsilon)) \},$$

and

$$Af(\eta^\epsilon(v)) = \frac{1}{\tau^\epsilon} \sum_{k=0}^{\lfloor n(v+\tau^\epsilon) \rfloor - \lfloor nv \rfloor - 1} \mathbb{E}_v^\epsilon \{ f(\eta^\epsilon(v) + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(\eta^\epsilon(v)) \}.$$

Without loss of generality, we illustrate the proof with  $f(y \mid Y_{t-1}, X_t, z_{t-d})$  is continuous in  $y$ , where the derivation for general case is analogous. Denote  $m^\epsilon = \lfloor n(v+\tau^\epsilon) \rfloor - \lfloor nv \rfloor$ , by Theorem 15.3 of [6], we have

$$\hat{A}^\epsilon f^\epsilon(v) = \sum_{k=0}^{m^\epsilon-1} \mathbb{E}_v^\epsilon \{ f(\eta^\epsilon(v) + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(\eta^\epsilon(v)) \} / \tau^\epsilon + o_p(1). \quad (7.47)$$

Meanwhile, for sufficiently large  $n$  and  $\epsilon = 1/n$ ,

$$\text{pr}(z_{t-d^0} \in (\theta_i^0, \theta_i^0 + \kappa_i \epsilon]) = \text{pr}(z_{t-d^0} = (\theta_i^0)^+ \kappa_i \epsilon (1+o(1))) = \frac{1}{n} \pi_z(\theta_i^0) \kappa_i (1+o(1)).$$

Thus, together with Assumptions 4a) and 5 and equation (48) in Appendix M, [7], for all fixed  $X$  with  $\rho$ -mixing coefficient sequence  $\{\rho(k)\}$ , there exists  $K^* \geq \mathbb{E}|f(X + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(X)|$  such that

$$\begin{aligned} & \frac{1}{\tau^\epsilon} \sum_{k=0}^{m^\epsilon-1} \mathbb{E} \left[ \left| \mathbb{E}_v^\epsilon \{ f(X + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(X) \} - \mathbb{E} \{ f(X + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(X) \} \right| \right] \\ & \leq \frac{1}{\tau^\epsilon} \left[ \mathbb{E} \{ f(X + \zeta_{k+\lfloor nv \rfloor + 1}^\epsilon) - f(X) \}^2 \right]^{1/2} \sum_{k=0}^{m^\epsilon-1} \rho(k+1) \\ & \leq \frac{1}{\tau^\epsilon} K^* \{ \text{pr}(Z_{t-d_0} \in (\theta_i^0, \theta_i^0 + \kappa_i \epsilon]) \}^{1/2} \sum_{k=0}^{m^\epsilon-1} \rho(k+1) \\ & = (1+o(1)) \frac{K^*}{\sqrt{n\tau^\epsilon}} (\pi_z(\theta_i^0) \kappa_i)^{1/2} \sum_{k=0}^{m^\epsilon-1} \rho(k+1) \rightarrow 0, \end{aligned}$$

by taking  $\tau^\epsilon = n^{-b}$  for some  $b \in (0, 1/2)$  and  $m^\epsilon \rightarrow \infty$ . Combining with (7.47), we have

$$\hat{A}^\epsilon f^\epsilon(v) = Af(\eta^\epsilon(v)) + o_p(1). \quad (7.48)$$

As  $p\text{-lim}_{\tau^\epsilon \rightarrow 0} f^\epsilon(v) = f(\eta^\epsilon(v))$ , by Theorem 1 of [29],  $\eta^\epsilon(v)$  converges weakly to  $\mathcal{C}(v)$ . On the other hand,  $\eta^\epsilon(v)$  is the uniquely solution such that  $f(\eta(t)) - \int_0^t Af(\eta(s))ds$  is a martingale for any  $f$  with continuous second-derivative and

a compact support ([57])). Choose  $v = 1$ , then  $\tilde{\ell}_{1,i}(\kappa_i)$  converges in distribution to  $\tilde{\ell}_{1,i}^*(\kappa_i)$  in the Skorohod metric. Analogously, we have the weak convergence of  $\tilde{\ell}_{2,i}(\kappa_i)$  to  $\tilde{\ell}_{2,i}^*(\kappa_i)$  in the Skorohod metric; by Cramer–Wold device, the weak convergence of  $\tilde{\ell}_i(\kappa_i)$  to  $\tilde{\ell}_i^*(\kappa_i)$  is obtained. Thus, Combining with Theorem 3.1 in [55],  $n(\hat{\theta}_i - \theta_i^0)$  converges weakly to  $M_i^-$ , where  $[M_i^-, M_i^+)$  is the unique random interval that minimizes  $\tilde{\ell}_i^*(\kappa_i)$ .

Next we consider the convergence of  $\hat{\Theta}$ . Analogously, we construct  $\tilde{\ell}(\kappa) = \sum_{i=1}^r \tilde{\ell}_i(\kappa_i)$ . Recall that for  $\tilde{\ell}_i^*(\kappa_i)$ , as suggested by (5.3) and (5.4), it has an intensity  $\pi_z(\theta_i^0)$  and jumps

$$\zeta_s(y_s^*; Y_{s-1}^*, X_s^*, \Psi_{i+1}^0, \Psi_i^0)I(\kappa_i \geq 0) + \zeta_s(y_s^*; Y_{s-1}^*, X_s^*, \Psi_i^0, \Psi_{i+1}^0)I(\kappa_i < 0),$$

with a non-degenerate distribution. Meanwhile, we have deduced the weak convergence of  $\tilde{\ell}_i(\kappa_i)$  to  $\tilde{\ell}_i^*(\kappa_i)$  from operator convergence as (7.48). Employing the idea in Theorem 3.3 of [32], for any constants  $c_1$  and  $c_2$  with either one nonzero, and any vectors  $\kappa^{(1)}, \dots, \kappa^{(4)} = \{\kappa_1^{(1)}, \dots, \kappa_r^{(1)}\}, \dots, \{\kappa_1^{(4)}, \dots, \kappa_r^{(4)}\}$  which are not all equal, the process

$$c_1[\tilde{\ell}(\kappa^{(1)}) - \tilde{\ell}(\kappa^{(2)})] + c_2[\tilde{\ell}(\kappa^{(3)}) - \tilde{\ell}(\kappa^{(4)})]$$

has a positive jump rate with bounded and non-degenerated jump sizes. Furthermore, operator convergence results analogous to (7.48) can be shown for this process, as both  $A$  and  $\hat{A}^\epsilon$  are linear operators. Thus, by similar arguments in the proof of Theorem 3.3 in [32], the weak convergence of  $\tilde{\ell}(\kappa)$  to the compound Poisson process  $\tilde{\ell}^*(\kappa)$  is established. Using Theorem 3.1 in [55] again,  $n(\hat{\Theta} - \Theta^0)$  converges weakly to  $M^-$ , where  $[M^-, M^+)$  is the unique minimizer of  $\tilde{\ell}^*(\kappa)$ .

In the general case, as  $\{\eta^\epsilon(v), 0 \leq v \leq 1\}$  is tight, every subsequence has a convergent subsequence, where the convergence can be assumed almost surely by enlarging the probability space. Moreover, as (7.47)–(7.48) hold, by using similar arguments as in [54],  $n(\hat{\Theta} - \Theta^0)$  converges weakly to  $M^-$ , where  $[M^-, M^+)$  is an almost surely minimizer of  $\tilde{\ell}^*(\kappa)$ .  $\square$

*Proof of Theorem 5.4.* For continuous  $z_t$ , the asymptotic independence between  $n(\hat{\Theta} - \Theta^0)$  and  $n^{1/2}(\hat{\Psi} - \Psi^0)$  follows from  $O_p(n^{-1})$  convergence rate of  $\hat{\Theta}$ , (7.45) and arguments of proving Theorem 2 of [10]. For discrete  $z_t$ , as we have obtained a very rapid convergence of  $\hat{\Theta}$  to  $\Theta^0$  in probability, the asymptotic independence between  $n(\hat{\Theta} - \Theta^0)$  and  $n^{1/2}(\hat{\Psi} - \Psi^0)$  can be similarly established.

Next we show the asymptotic distribution of model parameter estimates using techniques in [62]. Under Assumption 2 and the strong consistency of model parameters estimates, the estimator  $\hat{\Psi}$  satisfies

$$\begin{aligned} L'(\hat{\Psi}, \hat{\Theta}, d^0) &= \left( \frac{\partial L_{n,1}(\hat{\Psi}_1, \theta_0, \hat{\theta}_1, d^0)}{\partial \Psi_1}, \dots, \frac{\partial L_{n,r^0+1}(\hat{\Psi}_{r^0+1}, \hat{\theta}_{r^0}, \theta_{r^0+1}, d^0)}{\partial \Psi_{r^0+1}} \right)^\top \\ &= 0. \end{aligned}$$

By Taylor's expansion,

$$\begin{aligned} L'(\hat{\Psi}, \hat{\Theta}, d^0) &= L'(\Psi^0, \hat{\Theta}, d^0) + L''(\Psi^0, \hat{\Theta}, d^0)(\hat{\Psi} - \Psi^0) + o_p(1) \\ &= L'(\Psi^0, \Theta^0, d^0) + L''(\Psi^0, \Theta^0, d^0)(\hat{\Psi} - \Psi^0) + o_p(1), \end{aligned}$$

where the second-order partial derivative is symmetric. Next we verify the regularity conditions in Theorem 5.41 of [62]. First, from (4.1), for all  $i \neq j$  we have

$$\frac{\partial^2 L(\Psi, \Theta^0, d^0)}{\partial \Psi_i \partial \Psi_j} = 0, \quad \mathbb{E} \left[ \frac{\partial L(\Psi, \Theta^0, d^0)}{\partial \Psi_i} \left( \frac{\partial L(\Psi, \Theta^0, d^0)}{\partial \Psi_j} \right)^\top \right] = 0.$$

Thus,  $L''(\Psi^0, \Theta^0, d^0)$  and  $\mathbb{E} [L'(\Psi^0, \Theta^0, d^0)\{L'(\Psi^0, \Theta^0, d^0)\}^\top]$  are block-diagonal. In addition, by ergodic theorem,

$$\{L''(\Psi^0, \Theta^0, d^0)\}^{-1} \longrightarrow [\mathbb{E} \{L''(\Psi^0, \Theta^0, d^0)\}]^{-1},$$

in probability. From Assumption 2, we have the finiteness of  $\mathbb{E}[\{\partial L(\Psi^0, \Theta^0, d^0)/\partial \Psi_i\}^2]$  and  $\mathbb{E}[L'(\Psi^0, \Theta^0, d^0)\{L'(\Psi^0, \Theta^0, d^0)\}^\top]$ . Furthermore, Assumption 2 implies the boundedness of third-order derivatives of  $L(\Psi, \Theta^0, d^0)$  with respect to  $\Psi$  by some integrable function in the neighborhood of  $\Psi^0$ . Therefore, as all regularity conditions are satisfied, by Theorem 5.41 in [62], we have  $n^{1/2}(\hat{\Psi} - \Psi^0) \sim N(0, \Sigma^*)$ , where

$$\begin{aligned} \Sigma^* &= [\mathbb{E}\{L''(\Psi^0, \Theta^0, d^0)\}]^{-1} \mathbb{E} [L'(\Psi^0, \Theta^0, d^0)\{L'(\Psi^0, \Theta^0, d^0)\}^\top] \\ &\quad \times [\mathbb{E}\{L''(\Psi^0, \Theta^0, d^0)\}]^{-1}. \end{aligned}$$

As  $L'(\Psi^0, \Theta^0, d^0)$  is continuous in  $\Psi$  for sufficiently large  $n$ , by information matrix equality,

$$\mathbb{E} [L'(\Psi^0, \Theta^0, d^0)\{L'(\Psi^0, \Theta^0, d^0)\}^\top] = -[\mathbb{E}\{L''(\Psi^0, \Theta^0, d^0)\}]^{-1}.$$

Hence,  $\Sigma^*$  is defined as (5.5) and is diagonal. In particular,  $|\hat{\Psi} - \Psi^0| = O_p(n^{-1/2})$ .  $\square$

## References

- [1] ALBA, E. & TROYA, J. M. (1999). A survey of parallel distributed genetic algorithms. *Complexity* **4**, 31–52 [MR1688681](#)
- [2] ALBA, E. & TROYA, J. M. (2002). Improving flexibility and efficiency by adding parallelism to genetic algorithms. *Statistics and Complexity* **12**, 91–114 [MR1897509](#)
- [3] AN, H. Z. & HUANG, F. C. (1996). The geometrical ergodicity of nonlinear autoregressive models. *Statist. Sin.* **6**, 943–56 [MR1422412](#)
- [4] BRADLEY, R. C. (2005). Basic Properties of Strong Mixing Conditions. A Survey and Some Open Questions. *Probab. Surv.* **2**, 107–44 [MR2178042](#)

- [5] BILLINGSLEY, P. (1961). The Lindberg-Lévy theorem for martingales. *Proc. Amer. Math. Soc.* **12**, 788–92 [MR0126871](#)
- [6] BILLINGSLEY, P. (1968). *Convergence of Probability Measures*. New York: Wiley. [MR0233396](#)
- [7] BILLINGSLEY, P. (1999). *Convergence of Probability Measures* 2nd ed. New York: Wiley. [MR1700749](#)
- [8] CAI, Y. & STANDER, J. (2008). Quantile of self-excited threshold autoregressive time series models *J. Time Series Anal.* **29**, 186–202 [MR2387486](#)
- [9] CHAN, K. S. (1990). Deterministic stability, stochastic stability and ergodicity. In *Non-Linear Time Series: A Dynamical System Approach*, Ed. Tong, H. Oxford: Clarendon Press, Appendix 1 [MR1079320](#)
- [10] CHAN, K. S. (1993). Consistency and limiting distribution of the least square estimator of a threshold autoregressive model. *Ann. Statist.* **21**, 520–33 [MR1212191](#)
- [11] CHAN, K. S. & TONG, H. (1985). On the use of the deterministic Lyapunov function for the ergodicity of stochastic difference equations. *Adv. Appl. Prob.* **17**, 666–78 [MR0798881](#)
- [12] CHAN, K. S. & TONG, H. (1990). On likelihood ratio test for threshold autoregression. *J. R. S. S. B* **52**, 469–76 [MR1086798](#)
- [13] CHAN, K. S. & TSAY, R. S. (1998) Limiting properties of least square estimator of a continuous threshold model. *Biometrika* **85**, 413–26 [MR1649122](#)
- [14] CHAN, K. S. & GORACCI, G. (2019). On the ergodicity of first-order threshold autoregressive moving-average processes positive. *J. Time Series Anal.* **40**, 256–264 [MR3915530](#)
- [15] CHAN, N. H., YAU, C. Y. & ZHANG, R. M. (2015). LASSO estimation for threshold autoregressive models. *J. Econometrics* **189**, 285–96 [MR3414900](#)
- [16] CHAN, N. H., ING, C. K., LI, Y. & YAU, C. Y. (2017). Threshold estimation via group orthogonal greedy algorithm. *J. Bus. Econom. Stat.* **35**, 334–345 [MR3622841](#)
- [17] CRAMÉR, H. (1946). *Mathematical Methods of Statistics*. Princeton: Princeton University Press [MR0016588](#)
- [18] DAVIS, R. A., LEE, T. C. M. & RODRIGUEZ-YAM, G. A. (2006). Structural break estimation for nonstationary time series models. *J. Amer. Statist. Assoc.* **101**, 223–39 [MR2268041](#)
- [19] DAVIS, R. A. & YAU, C. Y. (2013). Consistency of minimum description length model selection for piecewise stationary time series models. *Electron. J. Stat.* **7**, 381–411 [MR3020426](#)
- [20] DAVIS, R. A., LEE, T. C. M. & RODRIGUEZ-YAM, G. A. (2008). Break detection for a class of nonlinear time series models. *J. Time Series Anal.* **29**, 834–67 [MR2450899](#)
- [21] DOUKHAN, P. (1994). *Mixing: Properties and Examples*. New York: Springer-Verlag [MR1312160](#)
- [22] FOKIANOS, K., RAHBEK, A. & TJØTHEIM, D. (2009). Poisson autoregression. *J. Amer. Statist. Assoc.* **104**, 1430–9 [MR2596998](#)
- [23] GAO, X., RITTER J. R. & ZHU, Z. (2013). Where have all IPO gone? *J. Fin. Quan. Anal.* **48**, 1663–92

- [24] IBRAGIMOV, I. A. (1963). A central limit theorem for a class of dependent random variables. *Theory of Probability & Its Application* **8**, 83–9 [MR0151997](#)
- [25] IBRANGIMOV, I. A. & KHASHMINSKII, R. Z. (1981). *Statistical Estimation: Asymptotic Theory*. New York: Springer-Verlag [MR0620321](#)
- [26] KONISHI, S. & KITAGAWA, G. (2008). *Information Criteria and Statistical Modeling*. New York: Springer Science+Business Media, LLC [MR2367855](#)
- [27] KOOP, G. & POTTER, S. M. (2003) Bayesian analysis of endogenous delay threshold models. *J. Bus.Econom. Stat.* **21**, 93–103 [MR1950381](#)
- [28] KURZ, T. G. (1975). Semigroups of conditioned shifts and approximation of Markov process. *Ann. Prob.* **3**, 618–42 [MR0383544](#)
- [29] KUSHNER, H. J. (1980). A martingale method for the convergence of a sequence of processes to a jump-diffusion process. *Z. Wahrscheinlichkeitstheorie verw. Gebiete* **53**, 207–19 [MR0580914](#)
- [30] LEE, T. C. M. (2000). An introduction to coding theory and the two-part minimum description length principle. *Int. Stat. Rev.* **69**, 169–83
- [31] LI, D., LI, W. K. & LING, S. (2011) On the least square estimation of threshold autoregressive and moving-average models. *Statistics and Its Inference* **4**, 183–196 [MR2812814](#)
- [32] LI, D. & LING, S. (2012). On the least square estimation of multiple-regime threshold autoregressive models. *J. Econometrics* **167**, 240–53 [MR2885449](#)
- [33] LI, D., LING, S. & ZAKOÏAN, J.-M. (2015). Asymptotic inference in multiple-threshold double autoregressive model. *J. Econometrics* **189**, 415–27 [MR3414910](#)
- [34] LI, D., LING, S. & ZHANG, R. (2016). On a Threshold Double Autoregressive Model. *J. Bus. Econom. Statist.* **34**, 161–75 [MR3450051](#)
- [35] LI, G. & LI, W. K. (2011) Testing a linear time series model against its threshold extension. *Biometrika* **98**, 243–50 [MR2804225](#)
- [36] LI, G., GUAN, B., LI, W. K. & YU, P. L. H. (2015) Hysteretic autoregressive time series model. *Biometrika* **102**, 717–23 [MR3394287](#)
- [37] LIU, J. & SUSKO, ED (1992) On strict stationarity and ergodicity of a non-linear ARMA model. *J. Appl. Prob* **29**, 363–373 [MR1165221](#)
- [38] LING, S. (1999). On the probabilistic properties of a double threshold ARMA conditional heteroskedastic model. *J. Appl. Prob.* **36**, 688–705 [MR1737046](#)
- [39] LOWRY, M. AND SCHWERT, G. W. (2002). IPO market cycles: bubbles or sequential earning? *J. Financ.* **3**, 1171–200
- [40] LOWRY, M. (2003). Why does IPO volume fluctuate so much? *J. Financ. Econom.* **67**, 3–40
- [41] LU, Q., LUND, R. & LEE, T. C. M. (2010). An MDL approach to the climate segmentation problem. *Ann. Appl. Statist.* **4**, 299–319 [MR2758173](#)
- [42] MCCULLAGH, P. & NELDER, J. A. (1989). *Generalized Linear Models*, 2nd Ed. London: Chapman and Hall [MR3223057](#)
- [43] MEYN S. P. & TWEEDIE R. L. (1993). *Markov Chains and Stochastic Stability*. London: Springer-Verlag. [MR1287609](#)
- [44] NEWEY, W. K. & MCFADDEN, D. (1994). Large Sample Estimation and



- Hypothesis Testing. In *Handbook of Econometrics, Vol.4*, Ed. Engle, R and McFadden, D. Elsevier Science. [MR1315971](#)
- [45] PÁSTOR, Ľ. & VERONESI, P. (2005). Rational IPO waves. *J. Financ.* **60**, 1713–57
- [46] PETROV, V. V.(1975). *Sums of Independent Random Variables*. Berlin: Springer-Verlag [MR0388499](#)
- [47] RISSANEN, J. (1989). *Stochastic Complexity in Statistical Inquiry*. Singapore: World Scientific [MR1082556](#)
- [48] RISSANEN, J. (2007). *Information and Complexity in Statistical Modeling*. Springer. [MR2287233](#)
- [49] RITTER, J. R. (1984). The “hot issue” market of 1980. *J. Bus.* **57**, 215–40
- [50] RITTER, J. R. & WELCH, I. (2002). A review of IPO activity, pricing and allocations. *J. Financ.* **4**, 1795–828
- [51] RITTER, J. R. (2013) Historical US IPO Statistics *Quandl.com* Retrieved June, 2015, from [https://www.quandl.com/data/RITTER/US\\_IPO\\_STATS-Historical-US-IPO-Statistics](https://www.quandl.com/data/RITTER/US_IPO_STATS-Historical-US-IPO-Statistics)
- [52] SAÏDI, Y. & ZAKOÏAN J.-M. (2006). Stationarity and Geometric Ergodicity of a Class of Nonlinear ARCH Models. *Ann. Appl. Probab.*, **4**, 2256–71 [MR2288721](#)
- [53] SAMIA, N. I., CHAN, K. S. & STENSETH, N. C. (2007). A generalized threshold mixed model for analyzing nonnormal nonlinear time series, with application to plague in Kazakhstan. *Biometrika* **94**, 101–18 [MR2367827](#)
- [54] SAMIA, N. I. & CHAN, K. S. (2011). Maximum likelihood estimation of a generalized threshold stochastic regression model. *Biometrika* **98**, 433–48 [MR2806439](#)
- [55] SEIJO, E. & SEN, B. (2011) A continuous mapping theorem for the smallest argmax functional. *Electron. J. Stat.* **5**, 421–39 [MR2802050](#)
- [56] SO, M. K. P., LI, W. K. & LAM, K. (2002). A Threshold Stochastic Volatility Model. *J. Forecast.* **21**, 473–500
- [57] STROOK, D. W. & VARADHAN, S. R. S. (1971). Diffusion process with boundary conditions. *Comm. Pure. Appl. Math.* **24**, 147–225 [MR0277037](#)
- [58] TONG, H. (1978). On a threshold model. In *Pattern Recognition and Signal Processing. NATO ASI series E: Applied Sc. (29)*. (ed. C.Chen), 575–86. Amsterdam: Sijthoff & Noordhoff
- [59] TONG, H. (1990). *Non-Linear Time Series. A Dynamical System Approach*. New York: Oxford University Press [MR1079320](#)
- [60] TONG, H. (2007). Birth of the threshold time series model. *Statist. Sinica.* **17**, 8–14
- [61] TONG, H. (2011). Threshold models in time series analysis-30 years on. *Stat. Interface* **4**, 107–18 [MR2812802](#)
- [62] VAN DER VAART, A. W. (2000). *Asymptotic Statistics*. Cambridge: Cambridge University Press [MR1652247](#)
- [63] WONG, C. S.& LI, W. K.(1997) Testing for threshold autoregression with conditional heteroskedasticity. *Biometrika* **84**, 407–18 [MR1467056](#)
- [64] WU, Y. & LING, S. (2017). Using genetic algorithms to parameters (d; r) estimation for threshold autoregressive models. *Computational Statistics*

- and Data Analysis* **35**, 318–333
- [65] YANG, Y. & LING, S. (2017). Inference for heavy-tailed and Multiple-Threshold Double Autoregressive Models. *J. Bus. Econom. Statist.* **35**, 318–333 [MR3622840](#)
- [66] YAU, C. Y., TANG, C. M. & LEE, T. C. M. (2015). Estimation of multiple-regime threshold autoregressive models with structural breaks. *J. Amer. Statist. Assoc.* **110**, 1175–86 [MR3420693](#)
- [67] YU, P. (2012). Likelihood estimation and inference in threshold regression. *J. Econometrics* **167**, 274–94 [MR2885451](#)
- [68] YUNG, C., ÇOLAK, G. & WANG, W. (2008). Cycles in the IPO market. *J. Financ. Econ.* **89**, 192–208