

Multiscale Autoregression on Adaptively Detected Timescales

Rafal Baranowski, Yining Chen and Piotr Fryzlewicz

London School of Economics and Political Science

Supplementary Materials

Contents

S1 Special cases of AMAR	2
S2 Additional numerical experiments	7
S3 Additional real data example: well-log	19
S4 Proofs of the theoretical results	23

S1 Special cases of AMAR

We now consider some special cases of AMAR, and offer visual insights into their behaviour.

S1.1 Special case I: a single scale

Let $\{X_t\}$ be a series following the AMAR model with a single scale, i.e.

$$X_t = \alpha_1 \frac{X_{t-1} + \dots + X_{t-\tau_1}}{\tau_1} + \varepsilon_t, \quad t = 1, \dots, T. \quad (\text{S1.1})$$

Recall that realisations for different values of α_1 (from 0.5 to 0.95, the latter corresponds to series that are near unit-root) and τ_1 (from 1 to 10) with standard Gaussian noise are plotted in Figure 1 of the main paper. It appears that the longer the scale, the noisier the appearance; the overall shape (driven by the low frequencies) is preserved, but the details (driven by the high frequencies) are increasingly obscured by noise. This behaviour can also be understood by considering the spectral properties of the single-scale AMAR model, where the fact that the corresponding AR coefficients in the single-scale AMAR model (S1.1) are constant provides a useful simplification in the form of the spectral density. With ε_t being white noise with unit variance, the spectral density of X_t given by (S1.1) is

$$S_X(f) = \left| 1 - \frac{\alpha_1}{\tau_1} e^{-2\pi f i} \frac{1 - e^{-2\pi f \tau_1 i}}{1 - e^{-2\pi f i}} \right|^{-2}, \quad |f| < \frac{1}{2}. \quad (\text{S1.2})$$

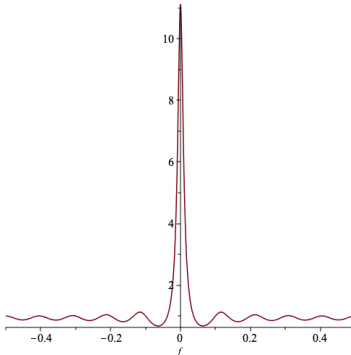


Figure 4: Spectral density of a single-scale AMAR process with $\tau_1 = 10$ and $\alpha_1 = 0.7$.

In view of the boundedness of $e^{-2\pi f i} \frac{1 - e^{-2\pi f \tau_1 i}}{1 - e^{-2\pi f i}}$ as a function of τ_1 , we have that $S_X(f) \rightarrow 1$ as $\tau_1 \rightarrow \infty$, for all $f \in (0, 1/2)$. However, as the sum of the AR coefficients in the single-scale AMAR does not depend on τ_1 , we have $S_X(0) = |1 - \alpha_1|^{-2}$. Note that more generally, given $\alpha_1, \dots, \alpha_q$, the spectral density at zero of any AMAR(q) process, i.e. its long-run variance, is independent of τ_1, \dots, τ_q .

As a visual illustration, Figure 4 shows the spectral density of a single-scale AMAR process with $\tau_1 = 10$ and $\alpha_1 = 0.7$. Due to the limiting behaviour described above, a single-scale AMAR for a large τ_1 can be approximated as the sum of two independent processes: one band-limited with a sharp peak at zero (and therefore representing a “slowly-varying” signal), and the other as white noise. This is in agreement with the appearance of the sample realisations shown in Figure 1, which begin to resemble a “signal + white noise” model for the larger values of τ_1 .

Finally, we note that even though all the series plotted in Figure 1 are weakly stationary, some of them exhibit behaviour that mimics non-stationarity, at least visually, when τ_1 is large, even for a moderate α_1 . This hints at the usefulness of AMAR in the

modelling of near unit-root or certain non-stationary series. More details can be found in a simulation study in Section S2.

S1.2 Special case II: one short plus one long scale

We now study the case of the AMAR model in which two timescales are present: one short one, and one long one. We have

$$X_t = \alpha_1 \frac{X_{t-1} + \dots + X_{t-\tau_1}}{\tau_1} + \alpha_2 \frac{X_{t-1} + \dots + X_{t-\tau_2}}{\tau_2} + \varepsilon_t. \quad (\text{S1.3})$$

First, if we keep $\alpha_1 + \alpha_2$ constant, and vary both coefficients from $\alpha_1 = 0$ on one extreme to $\alpha_2 = 0$ on the other extreme, then we obtain a “smooth transition” from a single-scale model with scale τ_2 to a single-scale model with scale τ_1 .

To gain further insight into the behaviour of AMAR with two timescales, now we consider $\alpha_1 = \alpha_2 = \alpha$, take $\tau_1 = 1$ and vary τ_2 . Figure 5 illustrates the case in which $\alpha_1 = \alpha_2 = \alpha = 0.49$, $\tau_1 = 1$ and $\tau_2 = 2, 10, 50$. When $\tau_2 = 50$, the longer scale has visually and practically no impact as the coefficients for the individual components (i.e. α_2/τ_2) are small. When $\tau_2 = 2$, we have a simple AR(2) model. On the other hand, when $\tau_2 = 10$, the realisation has the visual appearance of “a time-varying trend plus a low-order AR model”. Here, the longer scale is responsible for the changing “levels” at low-frequencies, whereas the shorter scale is responsible for instantaneous fluctuations at high-frequencies. This phenomenon is visually not present if $\sum_i \alpha_i$ is small or moderate, e.g. $\alpha = \alpha_1 = \alpha_2 = 0.3$, as demonstrated in Figure 5, but would show up if $\sum_i \alpha_i$ gets

moderately close to 1, e.g. at around 0.8.

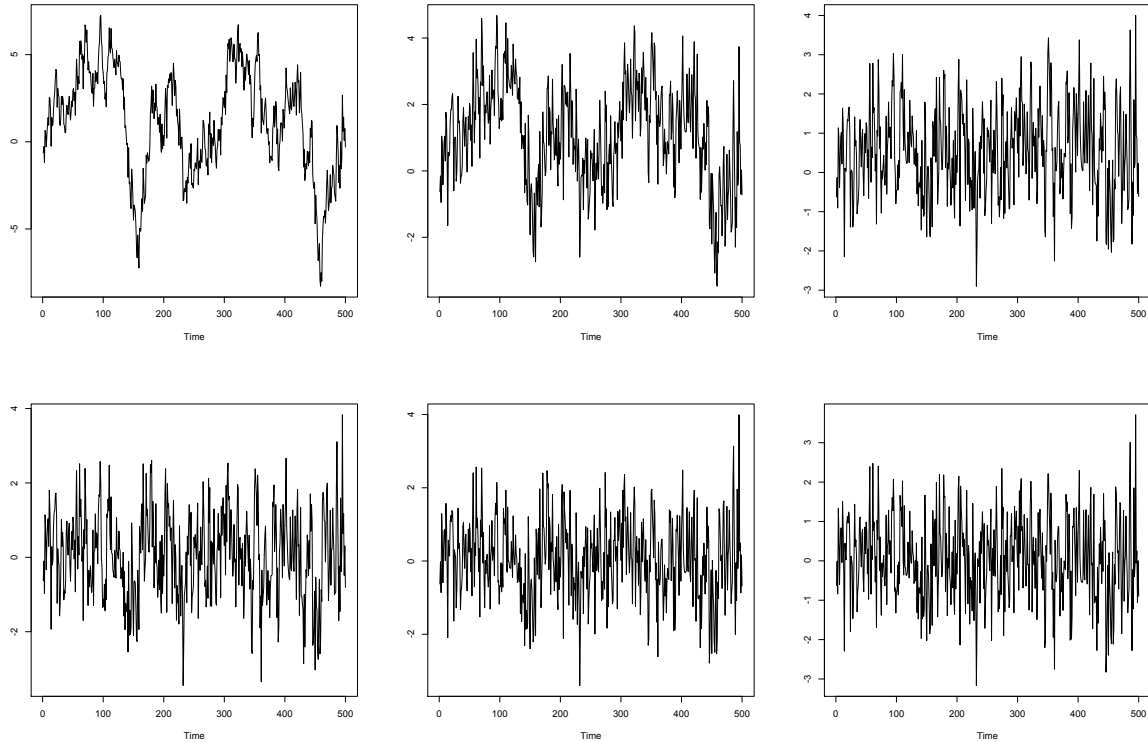


Figure 5: Simulated sample paths, of length 500, from the two-scale AMAR model (S1.3), with $\alpha_1 = \alpha_2 = \alpha = 0.49$ in the first row and $\alpha_1 = \alpha_2 = \alpha = 0.3$ in the second. Here $\tau_1 = 1$ and $\tau_2 = 2, 10, 50$ (respectively from left to right). The same random seed is used to generate path for each row.

Besides, models with a large τ_2 (for different α_i 's) also display the same interesting feature (i.e. “trend + noise” type, which might visually appear to be non-stationary), but it seems that the longer the scale τ_2 , the larger the sample sizes at which we are able to observe this phenomenon.

From these findings, we would infer that a two-scale AMAR model (with a small τ_1) is perhaps the most useful if (a) the longer scale is not too short or too long (i.e. in the order of 10s in practice), and (b) when the sum of the coefficients $\alpha_1 + \alpha_2$ is moderately

close to 1 (say > 0.8), with the coefficient α_2 of the longer scale not being too small. In this case the two-scale AMAR can imitate a time-varying trend plus a low order AR model, i.e. we are in a situation in which we are able to use a stationary AMAR model to model certain non-stationary-looking phenomena.

S1.3 Special case III: AMAR representation of seasonal models

In the class of seasonal $\text{ARIMA}(p, 0, q) \times (P, 0, Q)_S$ models, we consider models of the form $\Phi(B^S)\phi(B)X_t = \varepsilon_t$ (i.e. $q = Q = 0$), where

$$\begin{aligned}\Phi(B^S) &= 1 - \Phi_1 B^S - \dots - \Phi_P B^{PS} \\ \phi(B) &= 1 - \phi_1 B - \dots - \phi_p B^p,\end{aligned}$$

and where B is the lag operator. They belong to the class of AMAR models. As a simple example, consider $\text{ARIMA}(1, 0, 0) \times (1, 0, 0)_{12}$, an autoregressive model for monthly time series, with a single non-seasonal lag and yearly seasonality, given by

$$X_t = \phi_1 X_{t-1} + \Phi_1 X_{t-12} - \phi_1 \Phi_1 X_{t-13} + \varepsilon_t. \quad (\text{S1.4})$$

A typical characteristic feature of $\text{ARIMA}(p, 0, q) \times (P, 0, Q)_S$ models is its stretches of consecutive zero AR coefficients. For example, in (S1.4), the AR coefficients corresponding to lags 2 to 11 are zero. This means that AMAR models are also able to provide a relatively parsimonious representation of $\text{ARIMA}(p, 0, q) \times (P, 0, Q)_S$ models. As an

example, model (S1.4), represented in the AMAR framework, will need four scale parameters (at $\tau_1 = 1$, $\tau_2 = 11$, $\tau_3 = 12$ and $\tau_4 = 13$) and four corresponding AMAR coefficients (i.e. $\alpha_1 = \phi_1$, $\alpha_2 = -11\Phi_1$, $\alpha_3 = 12\Phi_1(1 + \phi_1)$ and $\alpha_4 = -13\phi_1\Phi_1$), which is more heavily parameterised than the *optimal* seasonal representation (S1.4) (with ϕ_1 and Φ_1) but much less than the full AR representation of (S1.4).

This (relative) parsimony of representation of seasonal models in the AMAR framework, plus the fact that the AMAR estimation framework is able to estimate the number of timescales and their spans automatically, makes AMAR a viable exploratory tool for identifying time series seasonality in data. In fact, we have demonstrated in the simulation study in Section 3.3 of the main paper that the AMAR estimation procedure is capable of identifying the right timescales rather effectively even with relatively small number of observations, confirming good potential of AMAR for the identification and exploratory analysis of seasonal models.

S2 Additional numerical experiments

S2.1 Sensitivity analysis

Several tuning parameters are required in the algorithm of our approach. The notable ones are the maximum number of scales q_{\max} , and the autoregressive order p used in the initial step. Besides, the choice of number of intervals M would also be required, but it should be apparent from our algorithm that it only plays a minor role under a large p (which, in the setup of our current algorithm, would imply $T > 250000$).

Based on our experiments, we find that the proposed approach is not too sensitive to the choice of all the aforementioned tuning parameters. Detailed results are given below.

S2.1.1 Maximum number of timescales – q_{\max}

Here we run the same experiments listed in the main manuscript, but set $q_{\max} = 5, 20$. The same evaluation metrics are used. Results are given in Table 6 and Table 7.

By comparing the results with those from Table 1 and Table 2 in the main manuscript (where by default $q_{\max} = 10$), it becomes evident that our approach does not appear to be sensitive to the choice of q_{\max} . In particular, for different choices of q_{\max} , every corresponding AMAR performs better than the competitors.

S2.1.2 The initial order of AR – p

We run the same experiments listed in the main manuscript, but use a fixed $p = 25$. The same evaluation metrics are used. Results are given in Table 8 and Table 9. For the ease of comparison, here we also recall the performance results of the default AMAR that uses p selected via SIC, for which details can be founded in Section 3.1 of the main manuscript.

Here we carefully fixed p at 25, so that it is larger than the timescales among all cases. Here the largest timescale is equal to $\lceil 3000^{0.4} \rceil = 24$, from Model (M6) with $T = 3000$. It can be seen that for most cases, both approaches perform similarly. Indeed, AMAR with a fixed p might lead to some very moderate improvement over our current approach of selection via SIC in a few settings. Still, not surprisingly, using a fixed p could be quite

Model (M1)								
q_{\max}	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	5	20	5	20	5	20	5	20
$T = 400$	0.164	0.17	0.539	0.589	0.016	0.0162	0.018	0.0134
	(0.013)	(0.013)	(0.043)	(0.046)	(0.00086)	(0.00082)	(0.0052)	(0.00092)
$T = 800$	0.051	0.051	0.187	0.206	0.00351	0.00385	0.00446	0.00469
	(0.0072)	(0.0074)	(0.032)	(0.034)	(0.00026)	(0.00036)	(0.00049)	(0.00054)
$T = 1500$	0.022	0.021	0.143	0.117	0.00116	0.00117	0.00138	0.00145
	(0.0046)	(0.0045)	(0.045)	(0.043)	(0.000088)	(0.000088)	(0.00024)	(0.00024)
$T = 3000$	0.01	0.011	0.021	0.049	0.000546	0.000549	0.000671	0.000685
	(0.0031)	(0.0033)	(0.0083)	(0.029)	(0.000027)	(0.000027)	(0.00017)	(0.00017)

Model (M2)								
q_{\max}	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	5	20	5	20	5	20	5	20
$T = 400$	0.251	0.289	1.07	1.27	0.0207	0.0206	0.0187	0.0243
	(0.017)	(0.017)	(0.064)	(0.071)	(0.0015)	(0.0014)	(0.002)	(0.0049)
$T = 800$	0.134	0.149	0.463	0.538	0.00551	0.00574	0.00649	0.00911
	(0.011)	(0.012)	(0.044)	(0.049)	(0.0006)	(0.00059)	(0.0011)	(0.0013)
$T = 1500$	0.125	0.136	1.18	1.26	0.00152	0.00142	0.00234	0.0025
	(0.011)	(0.011)	(0.13)	(0.13)	(0.00029)	(0.00026)	(0.00036)	(0.00039)
$T = 3000$	0.064	0.069	0.673	0.663	0.000159	0.000181	0.000776	0.00118
	(0.008)	(0.008)	(0.1)	(0.1)	(0.00011)	(0.000074)	(0.0002)	(0.00034)

Model (M3)								
q_{\max}	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	5	20	5	20	5	20	5	20
$T = 400$	0.511	0.706	1.46	1.35	0.0238	0.0209	0.0355	0.0289
	(0.023)	(0.035)	(0.054)	(0.046)	(0.00094)	(0.00075)	(0.0021)	(0.0016)
$T = 800$	0.262	0.344	0.631	0.64	0.00756	0.00701	0.0103	0.00918
	(0.018)	(0.026)	(0.034)	(0.034)	(0.00039)	(0.00031)	(0.00088)	(0.00075)
$T = 1500$	0.068	0.078	0.285	0.297	0.00197	0.00201	0.00341	0.00343
	(0.0089)	(0.011)	(0.04)	(0.042)	(0.0001)	(0.00011)	(0.00039)	(0.0004)
$T = 3000$	0.052	0.054	0.192	0.196	0.000677	0.000671	0.00152	0.00151
	(0.0078)	(0.0082)	(0.04)	(0.04)	(0.000042)	(0.000041)	(0.00023)	(0.00023)

Table 6: Performance of AMAR using different q_{\max} under (M1) – (M3), with estimated errors given in the brackets. Here \hat{q} is the number of the fitted timescales, D_H is the Hausdorff distance between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

Model (M4)								
q_{\max}	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	5	20	5	20	5	20	5	20
$T = 400$	0.065 (0.0079)	0.106 (0.013)	0.154 (0.022)	0.252 (0.031)	0.0104 (0.00085)	0.00932 (0.00061)	0.015 (0.0011)	0.0145 (0.001)
$T = 800$	0.041 (0.0063)	0.061 (0.0098)	0.129 (0.024)	0.131 (0.022)	0.00399 (0.0004)	0.00489 (0.00042)	0.00627 (0.00057)	0.00722 (0.00063)
$T = 1500$	0.041 (0.0063)	0.035 (0.006)	0.274 (0.053)	0.202 (0.046)	0.0018 (0.0001)	0.00193 (0.00025)	0.00313 (0.0004)	0.00352 (0.00056)
$T = 3000$	0.015 (0.0038)	0.023 (0.0049)	0.112 (0.037)	0.128 (0.036)	0.000753 (0.000023)	0.000752 (0.000023)	0.0017 (0.00024)	0.00173 (0.00024)
Model (M5)								
q_{\max}	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	5	20	5	20	5	20	5	20
$T = 400$	0.211 (0.017)	0.222 (0.018)	1.63 (0.072)	1.66 (0.073)	0.0107 (0.00043)	0.011 (0.00045)	0.0129 (0.00092)	0.0143 (0.0016)
$T = 800$	0.137 (0.013)	0.137 (0.013)	0.86 (0.055)	0.859 (0.055)	0.0042 (0.00023)	0.00421 (0.00022)	0.00535 (0.00056)	0.00533 (0.00056)
$T = 1500$	0.101 (0.012)	0.104 (0.013)	0.729 (0.078)	0.736 (0.078)	0.00167 (0.00011)	0.00171 (0.00012)	0.0023 (0.00034)	0.00223 (0.00034)
$T = 3000$	0.052 (0.0086)	0.052 (0.0086)	0.327 (0.054)	0.336 (0.056)	0.000343 (0.000044)	0.00034 (0.000043)	0.000757 (0.00018)	0.000761 (0.00018)
Model (M6)								
q_{\max}	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	5	20	5	20	5	20	5	20
$T = 400$	0.378 (0.022)	0.4 (0.023)	2.27 (0.054)	2.29 (0.053)	0.0133 (0.00048)	0.013 (0.00045)	0.0229 (0.0016)	0.0221 (0.0013)
$T = 800$	0.823 (0.031)	0.881 (0.035)	3.31 (0.072)	3.3 (0.071)	0.00909 (0.00029)	0.009 (0.00027)	0.0158 (0.00099)	0.0152 (0.00097)
$T = 1500$	0.428 (0.025)	0.462 (0.028)	3.09 (0.1)	3.1 (0.1)	0.00342 (0.00014)	0.00341 (0.00013)	0.00676 (0.00055)	0.00667 (0.00055)
$T = 3000$	0.533 (0.029)	0.644 (0.038)	3.57 (0.12)	3.52 (0.11)	0.00192 (0.000084)	0.00178 (0.000064)	0.00444 (0.00045)	0.00394 (0.00038)

Table 7: Performance of AMAR using different q_{\max} under (M4) – (M6), with estimated errors given in the brackets. Here \hat{q} is the number of the fitted timescales, D_H is the Hausdorff distance between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

problematic when the chosen p is close to or bigger than $\tau_{q_{\max}}$, as is evident in the setting of Model (M6) with $T = 3000$, where its performance is more than 100% worse in every evaluation metric.

S2.2 (More conventional) higher-order AR

In this part, we compare AMAR and conventional AR models (selected both by AIC and BIC) over the data that are generated from more conventional high-order stationary AR models. In particular, we consider the following settings, with $\tau_q = 16$ and $q = 16, 12, 8$.

(M7) $q = 16$ and $\tau_i = i$ for $i = 1, \dots, 16$, with the corresponding AR coefficients

$$\boldsymbol{\beta} = (0.2, -0.2, 0.2, -0.2, \dots, 0.2, -0.2)^T.$$

(M8) $q = 12$ and $\{\tau_1, \dots, \tau_{12}\} = \{1, \dots, 16\} \setminus \{2, 6, 10, 14\}$, with the corresponding AR coefficients

$$\boldsymbol{\beta} = (0.2, 0, 0, -0.2, 0.2, 0, 0, -0.2, \dots, 0.2, 0, 0, -0.2)^T.$$

(M9) $q = 8$ and $\tau_i = 2i$ for $i = 1, \dots, 8$, with the corresponding AR coefficients

$$\boldsymbol{\beta} = (0.2, 0.2, -0.2, -0.2, \dots, 0.2, 0.2, -0.2, -0.2)^T.$$

Here Model (M7) is a conventional high-order AR. Models (M8) and (M9) are also high-

Model (M1)								
p	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	SIC	fixed	SIC	fixed	SIC	fixed	SIC	fixed
$T = 400$	0.172 (0.014)	0.196 (0.016)	0.593 (0.047)	0.738 (0.07)	0.0159 (0.0008)	0.0178 (0.00091)	0.0133 (0.00093)	0.0147 (0.001)
$T = 800$	0.051 (0.0072)	0.047 (0.0071)	0.181 (0.03)	0.252 (0.046)	0.0035 (0.00026)	0.00401 (0.00032)	0.0046 (0.00048)	0.00483 (0.00052)
$T = 1500$	0.018 (0.0042)	0.02 (0.0046)	0.085 (0.03)	0.073 (0.027)	0.00116 (0.000088)	0.00115 (0.000084)	0.00138 (0.00024)	0.0014 (0.00025)
$T = 3000$	0.012 (0.0034)	0.009 (0.003)	0.072 (0.035)	0.016 (0.0063)	0.000546 (0.000027)	0.000546 (0.000027)	0.000662 (0.00017)	0.000681 (0.00017)
Model (M2)								
p	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	SIC	fixed	SIC	fixed	SIC	fixed	SIC	fixed
$T = 400$	0.303 (0.018)	0.235 (0.017)	1.33 (0.072)	1.67 (0.11)	0.02 (0.0013)	0.0233 (0.0018)	0.0281 (0.01)	0.0259 (0.007)
$T = 800$	0.194 (0.014)	0.154 (0.012)	0.764 (0.06)	1.13 (0.1)	0.00635 (0.00071)	0.00638 (0.00065)	0.00852 (0.0013)	0.00815 (0.0011)
$T = 1500$	0.108 (0.01)	0.122 (0.011)	0.921 (0.11)	0.821 (0.092)	0.00171 (0.00038)	0.000986 (0.00019)	0.00666 (0.0038)	0.00386 (0.0019)
$T = 3000$	0.07 (0.0081)	0.056 (0.0073)	0.646 (0.099)	0.446 (0.072)	0.0000979 (0.000021)	0.0000896 (0.000021)	0.000793 (0.0002)	0.000687 (0.00017)
Model (M3)								
p	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	SIC	fixed	SIC	fixed	SIC	fixed	SIC	fixed
$T = 400$	0.711 (0.035)	0.314 (0.024)	1.37 (0.046)	1.17 (0.057)	0.0211 (0.00076)	0.0187 (0.00072)	0.0296 (0.0016)	0.0297 (0.0017)
$T = 800$	0.344 (0.026)	0.146 (0.015)	0.643 (0.034)	0.481 (0.036)	0.00699 (0.00031)	0.00571 (0.00029)	0.00922 (0.00075)	0.00825 (0.00069)
$T = 1500$	0.083 (0.011)	0.087 (0.011)	0.31 (0.043)	0.254 (0.03)	0.00203 (0.00011)	0.0022 (0.00018)	0.0034 (0.0004)	0.00359 (0.00044)
$T = 3000$	0.054 (0.0082)	0.055 (0.0091)	0.219 (0.045)	0.126 (0.023)	0.000673 (0.000041)	0.000685 (0.000041)	0.0015 (0.00023)	0.00147 (0.00023)

Table 8: Performance of AMAR under (M1) – (M3) with the initial AR order p either selected via SIC, or fixed at $p = 25$. The estimated errors given in the brackets. Here \hat{q} is the number of the fitted timescales, D_H is the Hausdorff distance between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

Model (M4)								
p	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	SIC	fixed	SIC	fixed	SIC	fixed	SIC	fixed
$T = 400$	0.098 (0.012)	0.078 (0.012)	0.2 (0.027)	0.27 (0.041)	0.00892 (0.00065)	0.0105 (0.00074)	0.0145 (0.0011)	0.0157 (0.0012)
$T = 800$	0.044 (0.0085)	0.039 (0.008)	0.092 (0.019)	0.172 (0.035)	0.00397 (0.0003)	0.00446 (0.0004)	0.00657 (0.0006)	0.00653 (0.00057)
$T = 1500$	0.035 (0.006)	0.039 (0.0063)	0.291 (0.059)	0.158 (0.03)	0.00179 (0.00011)	0.00194 (0.00011)	0.00333 (0.0004)	0.00324 (0.00039)
$T = 3000$	0.023 (0.0051)	0.024 (0.005)	0.129 (0.033)	0.077 (0.019)	0.000756 (0.000023)	0.000753 (0.000023)	0.0017 (0.00024)	0.00162 (0.00024)
Model (M5)								
p	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	SIC	fixed	SIC	fixed	SIC	fixed	SIC	fixed
$T = 400$	0.217 (0.017)	0.265 (0.02)	1.64 (0.073)	2.29 (0.1)	0.0109 (0.00045)	0.0123 (0.00051)	0.0164 (0.0028)	0.0159 (0.0014)
$T = 800$	0.133 (0.013)	0.144 (0.014)	0.858 (0.056)	1.05 (0.073)	0.00414 (0.00022)	0.00447 (0.00026)	0.00517 (0.00055)	0.00588 (0.00065)
$T = 1500$	0.099 (0.012)	0.092 (0.011)	0.704 (0.076)	0.574 (0.058)	0.00167 (0.00012)	0.00168 (0.00011)	0.00237 (0.00033)	0.0023 (0.00034)
$T = 3000$	0.052 (0.0086)	0.049 (0.0082)	0.331 (0.054)	0.271 (0.042)	0.000339 (0.000043)	0.000346 (0.000043)	0.000788 (0.00017)	0.000746 (0.00017)
Model (M6)								
p	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $		$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$	
	SIC	fixed	SIC	fixed	SIC	fixed	SIC	fixed
$T = 400$	0.407 (0.024)	0.43 (0.024)	2.3 (0.054)	3.9 (0.12)	0.0133 (0.00046)	0.015 (0.00058)	0.023 (0.0016)	0.0279 (0.0016)
$T = 800$	0.886 (0.035)	0.486 (0.026)	3.29 (0.071)	3.18 (0.083)	0.00902 (0.00028)	0.00718 (0.00023)	0.015 (0.00098)	0.0129 (0.00086)
$T = 1500$	0.455 (0.028)	0.639 (0.035)	3.08 (0.1)	3.04 (0.085)	0.00336 (0.00013)	0.0038 (0.00014)	0.00668 (0.00055)	0.00712 (0.00056)
$T = 3000$	0.642 (0.037)	2.04 (0.063)	3.52 (0.11)	6.8 (0.12)	0.00177 (0.000064)	0.00392 (0.000096)	0.00395 (0.00038)	0.0071 (0.00055)

Table 9: Performance of AMAR under (M4) – (M6), with the initial AR order p either selected via SIC, or fixed at $p = 25$. Here \hat{q} is the number of the fitted timescales, D_H is the Hausdorff distance between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

order, but their AR coefficients are more structured (though $\tau_q = 16$ and q are still at the same order). In particular, all three models are stationary.

For these models, we run the experiments using the same settings as listed in the main manuscript, but set $q_{\max} = 20$ (as here q can be as high as 16). For the evaluation metrics, we look at the accuracy of the estimated order of AR, denoted by $|\hat{\tau}_q - \tau_q|$, the Euclidean distance between the fitted parameter vector and the true one, denoted by $\|\hat{\beta} - \beta\|$, and the mean squared prediction errors (MSPE) of different models. Results are given in Table 10.

We see that with in Model (M7), unsurprisingly AR with order selected via BIC performs the best among all the evaluation measures. However, the performance of AMAR is only slightly worse (and better than AR with order selected via AIC). In particular, it tends to estimates the number of scales (which is the same as the AR order) correctly when T is reasonably large, implying little efficiency loss for using AMAR even when there is no meaningful AMAR-type structure in the parameter vector of AR coefficients. On the other hand, as we move to Model (M8) and Model (M9) where the AR parameter vectors have more structures embedded (though here τ_q and q are still at the same order), AMAR tends to perform better than its competitors in terms of both the parameter estimation and prediction accuracy. The improvement is more visible in the setting of Model (M9), as it has less scales than Model (M8), so is intuitively more favourable to AMAR.

Model (M7)									
Method	$E \hat{\tau}_{\hat{q}} - \tau_q $			$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$		
	AMAR	AIC	BIC	AMAR	AIC	BIC	AMAR	AIC	BIC
$T = 400$	6.73 (0.14)	0.877 (0.042)	12.6 (0.17)	0.577 (0.0038)	0.051 (0.00084)	0.543 (0.0068)	0.204 (0.0031)	0.0534 (0.0018)	0.177 (0.003)
$T = 800$	2.05 (0.1)	1.68 (0.09)	0.23 (0.055)	0.139 (0.0053)	0.027 (0.00053)	0.0293 (0.0023)	0.0826 (0.0026)	0.0274 (0.0012)	0.0265 (0.0016)
$T = 1500$	0.916 (0.09)	2 (0.11)	0.014 (0.0037)	0.012 (0.00026)	0.0148 (0.0004)	0.0105 (0.00014)	0.0124 (0.00074)	0.0141 (0.00079)	0.0114 (0.00071)
$T = 3000$	0.662 (0.077)	2.14 (0.13)	0.015 (0.0038)	0.00582 (0.00015)	0.00733 (0.00017)	0.00532 (0.000074)	0.00597 (0.00052)	0.00701 (0.00057)	0.00555 (0.00049)
Model (M8)									
Method	$E \hat{\tau}_{\hat{q}} - \tau_q $			$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$		
	AMAR	AIC	BIC	AMAR	AIC	BIC	AMAR	AIC	BIC
$T = 400$	2.93 (0.085)	0.895 (0.04)	7.45 (0.21)	0.151 (0.0019)	0.0503 (0.00075)	0.175 (0.0038)	0.132 (0.003)	0.051 (0.0017)	0.14 (0.0039)
$T = 800$	1.62 (0.034)	1.69 (0.091)	0.307 (0.041)	0.08 (0.001)	0.0263 (0.00042)	0.025 (0.00077)	0.0767 (0.002)	0.0281 (0.0011)	0.0273 (0.0014)
$T = 1500$	0.669 (0.077)	2.02 (0.12)	0.019 (0.0052)	0.00969 (0.00029)	0.014 (0.00025)	0.0109 (0.00015)	0.00957 (0.0007)	0.0122 (0.00075)	0.0102 (0.00068)
$T = 3000$	0.597 (0.077)	1.95 (0.14)	0.014 (0.004)	0.00401 (0.00011)	0.00687 (0.00013)	0.00533 (0.000068)	0.00452 (0.00041)	0.0062 (0.0005)	0.00516 (0.00045)
Model (M9)									
Method	$E \hat{\tau}_{\hat{q}} - \tau_q $			$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$		
	AMAR	AIC	BIC	AMAR	AIC	BIC	AMAR	AIC	BIC
$T = 400$	1.37 (0.015)	0.849 (0.039)	0.161 (0.022)	0.104 (0.0011)	0.0499 (0.00082)	0.0485 (0.0012)	0.117 (0.003)	0.0493 (0.0017)	0.0487 (0.0017)
$T = 800$	1.04 (0.0061)	1.77 (0.093)	0.021 (0.0048)	0.0738 (0.00069)	0.0269 (0.00051)	0.0199 (0.00029)	0.0695 (0.002)	0.0253 (0.0011)	0.0212 (0.001)
$T = 1500$	0.214 (0.042)	1.93 (0.12)	0.017 (0.0041)	0.00469 (0.00039)	0.0147 (0.0003)	0.0107 (0.00016)	0.00691 (0.00067)	0.0131 (0.0008)	0.011 (0.00073)
$T = 3000$	0.174 (0.041)	1.75 (0.12)	0.018 (0.0047)	0.00215 (0.00022)	0.00703 (0.00015)	0.00528 (0.000079)	0.00384 (0.00044)	0.00646 (0.00051)	0.00551 (0.00047)

Table 10: Performance of different methods under (M7) – (M9), with estimated errors given in the brackets. Here $|\hat{\tau}_{\hat{q}} - \tau_q|$ is the difference between the estimated and true order of AR, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

S2.3 Non-stationary AR

Here we report the results from experiments with series simulated from non-stationary AR models with unit roots. The scenarios we consider are similar to (M1) – (M6) listed in the main manuscript, with their details outlined below.

(M1') Same as (M1) but with $\alpha_1 = 0.4$, $\alpha_2 = 0.6$ (i.e. $\boldsymbol{\beta} = (0.6, 0.2, 0.2)^T$).

(M2') Same as (M2) but with $\alpha_1 = 1.5$, $\alpha_2 = -0.5$ (i.e. $\boldsymbol{\beta} = (0.65, 0.65, -0.1, -0.1, -0.1)^T$).

(M3') Same as (M3) but with $\alpha_1 = 0.5$, $\alpha_2 = -1$, $\alpha_3 = 1.4$
(i.e. $\boldsymbol{\beta} = (0.5, -0.1, -0.1, -0.1, -0.1, 0.1, \dots, 0.1)^T$).

(M4') Same as (M4) but with $\alpha_1 = 1$, $\alpha_2 = -4.8$, $\alpha_3 = 10.2$, $\alpha_4 = -6.4$ (i.e. $\boldsymbol{\beta} = (1, 0, \dots, 0, 0.8, -0.8)^T$, so $\varepsilon_t = (1 - 0.8B^7)(1 - B)X_t$).

(M5') Same as (M5) but with $\alpha_1 = 1$ (i.e. $\boldsymbol{\beta} = (0.1, \dots, 0.1)^T$).

(M6') Same as (M6) but with $\alpha_1 = \alpha_2 = 0.5$ (i.e. $\boldsymbol{\beta} = (0.5+0.5/\lfloor T^{0.4} \rfloor, 0.5/\lfloor T^{0.4} \rfloor, \dots, 0.5/\lfloor T^{0.4} \rfloor)^T$).

Here we use AMAR with default choice of its tuning parameters outlined in Section 3.1. The corresponding results are summarised in Table 11 and Table 12, where as before, we report the estimates for $|q - \hat{q}|$, with \hat{q} being the number of the fitted timescales, the Hausdorff distance D_H between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, the Euclidean distance between the fitted parameter vector and the true one, denoted by $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|$, and the ratio between the mean squared prediction error (MPSE) using the fitted model and that with the oracle over the next $T^* = 100$ unseen observations.

Model (M1')										
Method	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$		
	AMAR	Fused	AMAR	Fused	AMAR	Fused	AIC	AMAR	Fused	AIC
$T = 400$	0.469	1.21	2.31	18.6	0.0449	0.381	0.0268	38.6	28.7	32.5
	(0.019)	(0.049)	(0.11)	(0.042)	(0.0022)	(0.0014)	(0.0008)	(23)	(1.2)	(3.1)
$T = 800$	0.33	1.43	1.97	26.4	0.0347	0.398	0.0201	1.45	29	7.79
	(0.016)	(0.081)	(0.11)	(0.063)	(0.0022)	(0.001)	(0.00068)	(1.4)	(1.3)	(0.6)
$T = 1500$	0.367	2.13	4.47	36.1	0.0302	0.409	0.017	0.0231	26.9	1.95
	(0.016)	(0.14)	(0.25)	(0.088)	(0.0022)	(0.00073)	(0.00059)	(0.0018)	(1.1)	(0.16)
$T = 3000$	0.249	2.61	3.44	51.8	0.0192	0.42	0.0161	0.0148	29.8	0.536
	(0.014)	(0.18)	(0.23)	(0.094)	(0.0019)	(0.00026)	(0.00063)	(0.0016)	(1.3)	(0.043)

Model (M2')										
Method	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$		
	AMAR	Fused	AMAR	Fused	AMAR	Fused	AIC	AMAR	Fused	AIC
$T = 400$	0.399	3.85	2.09	13.5	0.0254	0.609	0.116	0.0374	4.41	0.226
	(0.019)	(0.1)	(0.075)	(0.12)	(0.00093)	(0.0092)	(0.0031)	(0.0024)	(0.14)	(0.014)
$T = 800$	0.269	6.75	1.52	18.5	0.0134	0.69	0.0892	0.0175	6	0.0971
	(0.017)	(0.16)	(0.073)	(0.21)	(0.00061)	(0.0073)	(0.0025)	(0.0015)	(0.2)	(0.0055)
$T = 1500$	0.187	10	2.11	24.5	0.00572	0.733	0.0681	0.00601	7.33	0.0506
	(0.013)	(0.22)	(0.16)	(0.33)	(0.00036)	(0.0059)	(0.0021)	(0.00074)	(0.25)	(0.0024)
$T = 3000$	0.086	14.2	0.985	35.3	0.0018	0.774	0.0399	0.0023	8.7	0.0246
	(0.009)	(0.31)	(0.12)	(0.5)	(0.0002)	(0.0039)	(0.0014)	(0.00051)	(0.25)	(0.0013)

Model (M3')										
Method	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$		
	AMAR	Fused	AMAR	Fused	AMAR	Fused	AIC	AMAR	Fused	AIC
$T = 400$	0.747	2.47	1.37	17.4	0.0246	0.297	0.0606	0.0351	0.766	17.4
	(0.034)	(0.058)	(0.043)	(0.13)	(0.0011)	(0.003)	(0.00089)	(0.0021)	(0.022)	(2.1)
$T = 800$	0.499	1.72	0.897	24.7	0.0201	0.318	0.0342	0.0365	0.767	3.76
	(0.029)	(0.089)	(0.039)	(0.11)	(0.0015)	(0.0027)	(0.00059)	(0.0036)	(0.018)	(0.34)
$T = 1500$	0.177	2.4	0.744	34.1	0.0047	0.323	0.0216	0.0104	0.788	0.954
	(0.015)	(0.16)	(0.073)	(0.13)	(0.00076)	(0.0031)	(0.00043)	(0.0019)	(0.019)	(0.09)
$T = 3000$	0.104	1.93	0.506	49.2	0.00336	0.339	0.0141	0.00609	0.854	0.288
	(0.011)	(0.16)	(0.069)	(0.15)	(0.00073)	(0.0027)	(0.00035)	(0.0015)	(0.02)	(0.028)

Table 11: Performance of different methods under (M1') – (M3'), with estimated errors given in the brackets. Here \hat{q} is the number of the fitted timescales, D_H is the Hausdorff distance between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

Model (M4')											
Method	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$			
	AMAR	Fused	AMAR	Fused	AMAR	Fused	AIC	AMAR	Fused	AIC	
$T = 400$	0.527	2.64	1.45	18.1	0.161	2.22	0.888	0.254	216	28.1	
	(0.021)	(0.024)	(0.07)	(0.074)	(0.013)	(0.0019)	(0.013)	(0.021)	(9.3)	(2)	
$T = 800$	0.27	2.55	0.702	25.4	0.21	2.24	0.84	0.295	227	7.28	
	(0.017)	(0.026)	(0.052)	(0.11)	(0.014)	(0.00068)	(0.012)	(0.022)	(10)	(0.5)	
$T = 1500$	0.225	2.4	2.02	34.2	0.124	2.25	0.847	0.189	267	2.79	
	(0.014)	(0.041)	(0.15)	(0.17)	(0.011)	(0.00035)	(0.011)	(0.018)	(13)	(0.13)	
$T = 3000$	0.175	2.97	1.69	48.2	0.0678	2.26	0.841	0.0921	293	1.4	
	(0.012)	(0.092)	(0.14)	(0.22)	(0.0085)	(0.00026)	(0.0096)	(0.012)	(14)	(0.052)	
Model (M5')											
Method	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$			
	AMAR	Fused	AMAR	Fused	AMAR	Fused	AIC	AMAR	Fused	AIC	
$T = 400$	0.354	0.575	1.94	9.74	0.0134	0.0461	0.0464	5.81	0.439	10	
	(0.022)	(0.037)	(0.077)	(0.039)	(0.00061)	(0.00035)	(0.00082)	(4.5)	(0.032)	(0.89)	
$T = 800$	0.322	1.33	1.44	17.5	0.00602	0.0591	0.0276	0.254	0.467	2.47	
	(0.02)	(0.065)	(0.074)	(0.07)	(0.00052)	(0.00046)	(0.00062)	(0.19)	(0.027)	(0.19)	
$T = 1500$	0.154	2.3	1.2	27.6	0.00212	0.069	0.0203	2.32	0.569	0.692	
	(0.013)	(0.12)	(0.11)	(0.062)	(0.0003)	(0.00044)	(0.0006)	(2.2)	(0.034)	(0.064)	
$T = 3000$	0.114	4.19	0.818	43.4	0.000744	0.0776	0.014	0.72	0.622	0.197	
	(0.01)	(0.19)	(0.088)	(0.1)	(0.00024)	(0.00036)	(0.00048)	(0.54)	(0.03)	(0.022)	
Model (M6')											
Method	$E \hat{q} - q $		$E(D_H)$		$E\ \hat{\beta} - \beta\ $			$\frac{\text{MSPE}(\text{fitted})}{\text{MSPE}(\text{oracle})} - 1$			
	AMAR	Fused	AMAR	Fused	AMAR	Fused	AIC	AMAR	Fused	AIC	
$T = 400$	0.889	1.32	3.67	17.5	0.0177	0.234	0.0441	0.0251	0.762	10.4	
	(0.034)	(0.064)	(0.076)	(0.11)	(0.0007)	(0.002)	(0.00068)	(0.0013)	(0.038)	(0.87)	
$T = 800$	1.54	1.52	7.63	24.2	0.0184	0.228	0.0329	0.0336	0.631	3.37	
	(0.049)	(0.11)	(0.072)	(0.17)	(0.00073)	(0.0022)	(0.00041)	(0.0015)	(0.031)	(0.3)	
$T = 1500$	0.931	2.01	5.6	33.4	0.00501	0.23	0.0229	0.01	0.499	1.07	
	(0.04)	(0.15)	(0.13)	(0.22)	(0.00025)	(0.0022)	(0.00025)	(0.00075)	(0.018)	(0.094)	
$T = 3000$	2.09	2.95	12.4	47.9	0.00515	0.236	0.0159	0.0122	0.427	0.328	
	(0.066)	(0.18)	(0.13)	(0.28)	(0.00011)	(0.0019)	(0.00015)	(0.00076)	(0.012)	(0.036)	

Table 12: Performance of different methods under (M4') – (M6'), with estimated errors given in the brackets. Here \hat{q} is the number of the fitted timescales, D_H is the Hausdorff distance between the fitted timescale locations $\{\hat{\tau}_1, \dots, \hat{\tau}_{\hat{q}}\}$ and the true ones $\{\tau_1, \dots, \tau_q\}$, $\|\hat{\beta} - \beta\|$ is the Euclidean distance between the fitted parameter vector and the true one, and MPSE is the mean squared prediction errors of different models.

We see that even in the setting of non-stationary observations, AMAR still performs much better than its competitors in most settings, even though all methods seem to perform worse as compared to the stationary settings. Unsurprisingly, here the reported results are associated with larger estimation errors.

In addition, we note that the fused LASSO approach performs much worse than its competitors in terms of MSPE, especially in (M1') and (M4'). This is because the fused LASSO approach tends to over-estimate the number of scales, resulting in less accurate $\hat{\beta}$, which could greatly affect the corresponding MSPE when the series is non-stationary.

S3 Additional real data example: well-log

We consider the well-log data from O Ruanaidh and Fitzgerald (1996). Prior to use, the data is cleaned by removing outliers, taken here to be the observations that differ from the median-filter fit to the data (with span 25) by at least 7500. This retains 97.7% of the data points. The cleaned data, denoted as $\{X_t\}_{t=1}^{3956}$, is shown in the left plot of Figure 6.

As summarised in Fearnhead and Clifford (2003), the data represents measurements of the nuclear magnetic response of underground rocks. The underlying (unobserved) signal is assumed to be piecewise constant, with each constant segment representing a stratum of rock. The jumps occur when a new rock stratum is met. The problem of detecting these change-points in the underlying signal is of practical importance in oil drilling.

It is known (for instance, see Cho and Fryzlewicz (2021) and the references therein)

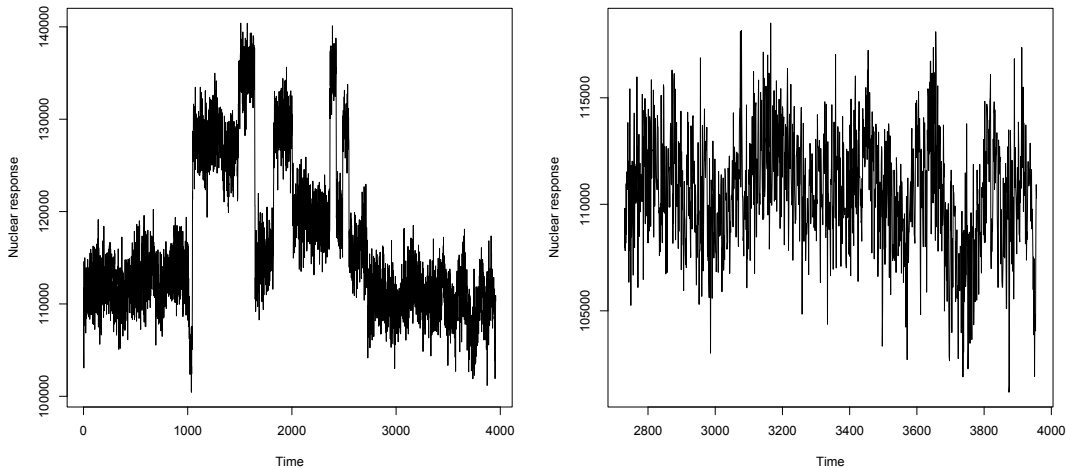


Figure 6: Left: the well-log data from O Ruanaidh and Fitzgerald (1996), cleaned as described in the text. Right: the end part of the data, from time location 2730.

that the problem of multiple change-point detection in a piecewise-constant signal observed in noise is much more challenging if the noise displays autocorrelation, as the natural fluctuations of the autocorrelated process can be mistaken for change-points, and vice versa. This appears to be the case in the well-log data: the right-hand plot of Figure 6 shows the end portion of the data, from the observation after the last visually obvious change-point (at location 2729) to the end. As discussed earlier, the visual appearance of the data fluctuations in this region of the dataset suggests that the AMAR model may be appropriate. Our aim is therefore to: (a) estimate the appropriate AMAR model on $\{X_{2730}, \dots, X_{3956}\}$, (b) fit the estimated model from the previous step on the entire dataset (i.e. $\{X_1, \dots, X_{3956}\}$) to remove the autocorrelations in the data, and (c) estimate change-point locations in the thus-decorrelated dataset using a method suitable for multiple change-point detection in uncorrelated (Gaussian) noise.

We start with a preliminary time series analysis of $\{X_{2730}, \dots, X_{3956}\}$. The unconstrained AR fit to this subset of the data, with the AR order chosen via AIC yields order 17, and the estimated coefficients are shown in the left panel of Figure 7. The appearance of the vector of the estimated coefficients suggests that a piecewise-constant model (as dictated by AMAR) may be suitable here. The fitted AMAR model returns estimated scales 1, 9, 13, 16, 17 (see Figure 7).

Prior to fitting the estimated AMAR model to the entire dataset, however, we shrink the estimated AMAR coefficients by a factor of $\rho \in (0, 1)$, i.e. we replace each estimated AMAR coefficient $\hat{\alpha}_r$ by $\rho\hat{\alpha}_r$. This is done because the original estimated AMAR coefficients sum up to practically 1 (0.9998), and therefore fitting such a “near-unit-root” AMAR model has a strong differencing effect, which as well as successfully removing the autocorrelations, could also potentially remove too much of the structure of the signal for successful detection of change-points in the levels.

We choose ρ as follows. Starting with $\rho = 0$, we increase ρ in steps of 0.01 until our selected procedure(s) for change-point detection under lack of serial correlation do not indicate any change-points after time $t = 2730$ (since we initially fitted an AMAR model on this portion of the data under the assumption of stationarity there). This is first achieved for $\rho = 0.78$, for both Wild Binary Segmentation (Fryzlewicz, 2014) and Narrowest-Over-Threshold (Baranowski et al., 2019), both with model selection via the strengthened Schwarz Information Criterion, and using the implementation from the R package **breakfast** (Anastasiou et al., 2021) with otherwise default parameters. These two procedures indicate, respectively, 12 and 10 change-points in the signal. The change-

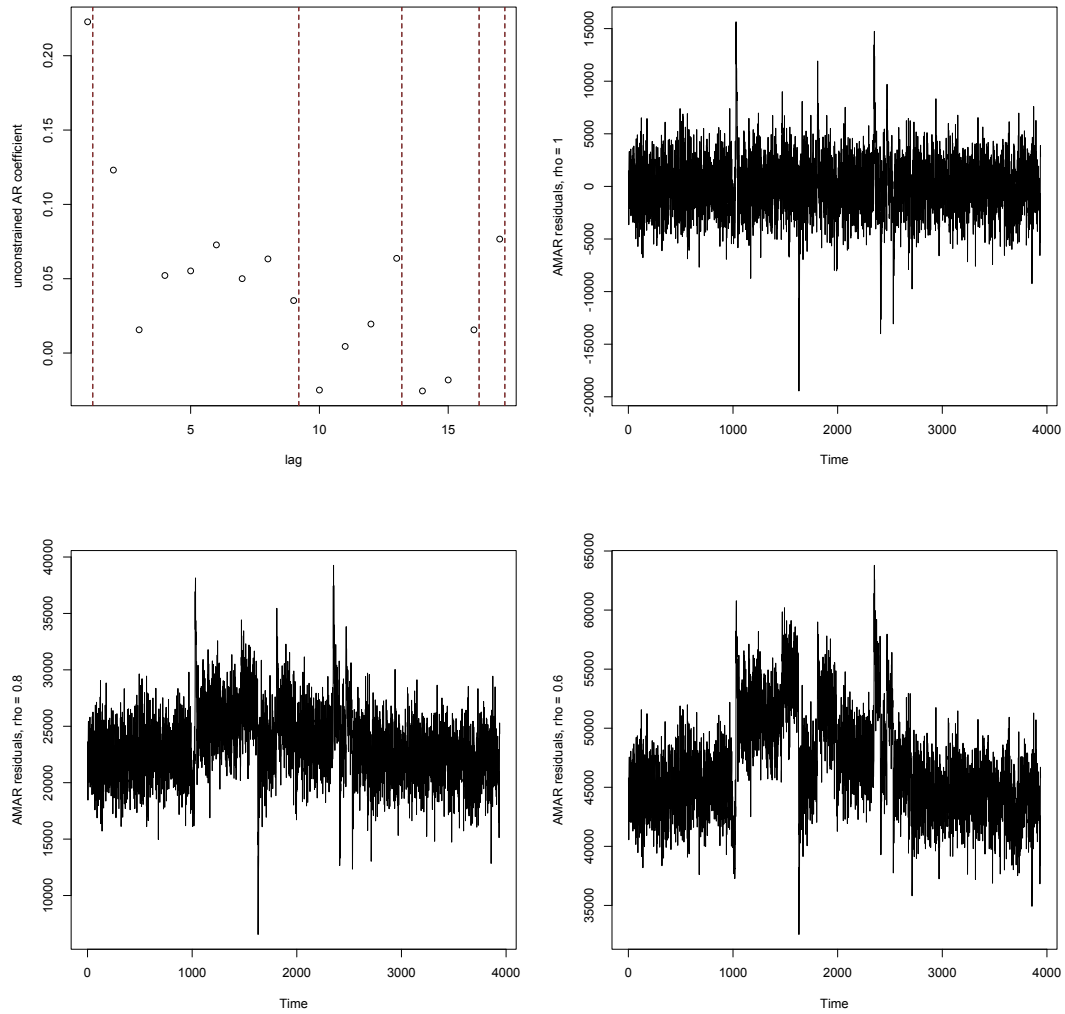


Figure 7: Top left: the unconstrained estimated AR coefficients for $X_{2730:3956}$ (circles); extents of estimated AMAR scales (dashed lines). Top right: unshrunk AMAR residuals. Bottom left: AMAR residuals with $\rho = 0.8$. Bottom right: AMAR residuals with $\rho = 0.6$.

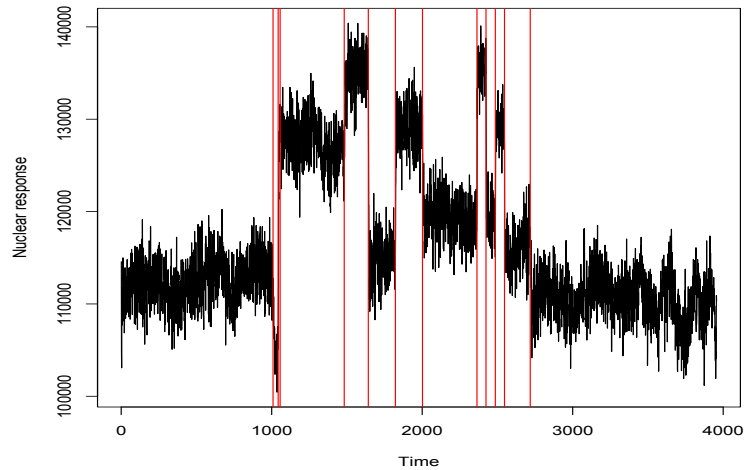


Figure 8: The well-log data with the change-point locations estimated via the shrunk AMAR fit with $\rho = 0.78$ and using WBS+sSIC on the residuals.

point locations estimated via Wild Binary Segmentation are shown in Figure 8. With the exception of the possible double detection at times $t = 1043, 1056$, the estimated change-point locations visually align with the signal very well.

S4 Proofs of the theoretical results

S4.1 Proof of Proposition 2.1

For $\text{AR}(p)$ processes, it has a stationary and causal solution if and only if all the roots of $b(z) = 0$ lie outside \mathbb{T} .

For any $\text{AMAR}(q)$ with $\alpha_1, \dots, \alpha_q$ (and the corresponding AR parameters β_1, \dots, β_p),

$\sum_{j=1}^q |\alpha_j| < 1$ under the AMAR framework is equivalent to

$$\sum_{j=1}^p |\beta_j| = \sum_{j=1}^p \left(\sum_{k:\tau_k \geq j} \frac{|\alpha_k|}{\tau_k} \right) = \frac{|\alpha_1|}{\tau_1} \tau_1 + \cdots + \frac{|\alpha_q|}{\tau_q} \tau_q < 1$$

in view of Equations (1.1) and (2.2) and (2.3). Now since $\sum_{j=1}^p |\beta_j| < 1$, $b(z) := 1 - \beta_1 z - \cdots - \beta_p z^p \geq 1 - |\beta_1| \|z\| - \cdots - |\beta_p| \|z\|^p \geq 1 - \sum_{j=1}^p |\beta_j| > 0$ for any $z \in \mathbb{T}$. As such, all the roots of $b(z) = 0$ lie outside \mathbb{T} , which implies the existence of a causal stationary solution.

Next, given $\alpha_1, \dots, \alpha_q \geq 0$, we have that $\beta_1, \dots, \beta_p \geq 0$. The existence of a causal stationary solution implies that all the roots of $b(z) = 0$ lie outside \mathbb{T} . Since $b(0) = 1$ and $b(\cdot)$ is continuous, one would necessarily require $b(1) > 0$. i.e. $\beta_1 + \cdots + \beta_p < 1$. This condition under the AMAR framework is equivalent to

$$\sum_{j=1}^p \left(\sum_{k:\tau_k \geq j} \frac{\alpha_k}{\tau_k} \right) = \frac{\alpha_1}{\tau_1} \tau_1 + \cdots + \frac{\alpha_q}{\tau_q} \tau_q < 1,$$

which is the same as $\sum_{j=1}^q |\alpha_j| < 1$ under non-negativity. \square

S4.2 Proof of Theorem 2.1

We write the AR(p) model as

$$\mathbf{Y}_t = \mathbf{B}\mathbf{Y}_{t-1} + \varepsilon_t \mathbf{u}, \quad t = 1, \dots, T, \quad (\text{S4.5})$$

where $\mathbf{Y}_t = (X_t, X_{t-1}, \dots, X_{t-p+1})'$, the matrix of the coefficients

$$\mathbf{B} = \begin{pmatrix} \beta_1 & \beta_2 & \cdots & \beta_p \\ & \mathbf{I}_{p-1} & & \mathbf{0} \end{pmatrix} \quad (\text{S4.6})$$

and $\mathbf{u} = (1, 0, \dots, 0)' \in \mathbb{R}^p$. We start with a few auxiliary results.

Lemma S4.1 (Parseval's identity, Theorem 1.9 in Duoandikoetxea (2001)) *For any complex-valued sequence $\{f_k\}_{k \in \mathbb{Z}}$ such that $\sum_{k \in \mathbb{Z}} |f_k|^2 < \infty$, the following identity holds*

$$\sum_{k \in \mathbb{Z}} |f_k|^2 = \int_{\mathbb{T}} |f(z)|^2 dm(z), \quad (\text{S4.7})$$

where $f(z) = \sum_{k \in \mathbb{Z}} f_k z^k$, $\mathbb{T} = \{z \in \mathbb{C} : |z| = 1\}$, $dm(z) = \frac{d|z|}{2\pi}$.

Lemma S4.2 (Cauchy's integral formula) *Let $\mathbf{M} \in \mathbb{R}^{p \times p}$ be a real- or complex- valued matrix. Then for any curve Γ enclosing all eigenvalues of \mathbf{M} and any $j \in \mathbb{N}$ the following holds*

$$\mathbf{M}^j = \frac{1}{2\pi i} \int_{\Gamma} z^j (z\mathbf{I}_p - \mathbf{M})^{-1} dz = \frac{1}{2\pi i} \int_{\Gamma} z^{j-1} (\mathbf{I}_p - z^{-1}\mathbf{M})^{-1} dz. \quad (\text{S4.8})$$

Lemma S4.3 *Let \mathbf{B} given by (S4.6) be the matrix of coefficients of a stationary AR(p) process and let $\mathbf{v} = (v_1, \dots, v_p)' \in \mathbb{R}^p$. For all $z \in \mathbb{C}$ such that $\sum_{i=0}^{\infty} |\langle \mathbf{v}, \mathbf{B}^i \mathbf{u} \rangle| |z^i| < \infty$,*

we have

$$b(z) \sum_{i=0}^{\infty} \langle \mathbf{v}, \mathbf{B}^i \mathbf{u} \rangle z^i = b(z) \langle \mathbf{v}, (\mathbf{I}_p - z\mathbf{B})^{-1} \mathbf{u} \rangle = v(z), \quad (\text{S4.9})$$

where $v(z) = v_1 + v_2 z + \dots + v_p z^{p-1}$, and where $b(z)$ is the AR polynomial.

Proof. As $\sum_{i=0}^{\infty} |\langle \mathbf{v}, \mathbf{B}^i \mathbf{u} \rangle| |z|^i < \infty$, we can change the order of summation in the left-hand side of (S4.9)

$$(1 - \beta_1 z - \dots - \beta_p z^p) \sum_{i=0}^{\infty} \langle \mathbf{v}, \mathbf{B}^i \mathbf{u} \rangle z^i = \left\langle \mathbf{v}, \left(\sum_{i=0}^{\infty} (1 - \beta_1 z - \dots - \beta_p z^p) z^i \mathbf{B}^i \right) \mathbf{u} \right\rangle.$$

Define $\beta_0 = -1$, $\beta_k = 0$ for $k > p$. By direct algebraic manipulation,

$$\sum_{i=0}^{\infty} (1 - \beta_1 z - \dots - \beta_p z^p) z^i \mathbf{B}^i = - \sum_{i=0}^{\infty} \left(\sum_{k=0}^i \beta_k \mathbf{B}^{i-k} \right) z^i := - \sum_{i=0}^{\infty} \mathbf{D}_i z^i.$$

The characteristic polynomial of \mathbf{B} is given by $\phi(z) = \sum_{k=0}^p \beta_k z^{p-k}$. From the Cayley–Hamilton theorem, \mathbf{B} is a root of ϕ , and, consequently for $i \geq p$,

$$\mathbf{D}_i = \mathbf{B}^{i-p} \sum_{k=0}^i \beta_k \mathbf{B}^{p-k} = \mathbf{B}^{i-p} \sum_{k=0}^p \beta_k \mathbf{B}^{p-k} = 0.$$

It remains to demonstrate that $\langle \mathbf{v}, \mathbf{D}_i \mathbf{u} \rangle = -v_{i+1}$ for $i = 0, \dots, p-1$, which we show by induction. For $i = 0$, $\langle \mathbf{v}, \mathbf{D}_0 \mathbf{u} \rangle = \beta_0 \langle \mathbf{v}, \mathbf{u} \rangle = -v_1$. When $i \geq 1$, matrices \mathbf{D}_i satisfy

$\mathbf{D}_i = \mathbf{B}\mathbf{D}_{i-1} + \beta_i \mathbf{I}_p$, therefore

$$\begin{aligned} \langle \mathbf{v}, \mathbf{D}_i \mathbf{u} \rangle &= \langle \mathbf{v}, \mathbf{B}\mathbf{D}_{i-1} \mathbf{u} \rangle + \beta_i \langle \mathbf{v}, \mathbf{u} \rangle = \langle \mathbf{B}' \mathbf{v}, \mathbf{D}_{i-1} \mathbf{u} \rangle + \beta_i \langle \mathbf{v}, \mathbf{u} \rangle \\ &= \langle v_1(\beta_1, \dots, \beta_p)' + (0, v_2, \dots, v_p)', \mathbf{D}_{i-1} \mathbf{u} \rangle + \beta_i \langle \mathbf{v}, \mathbf{u} \rangle = -v_1 \beta_i - v_{i+1} + v_1 \beta_i \\ &= -v_{i+1}, \end{aligned}$$

which completes the proof. \square

Lemma S4.4 *Let Z_1, Z_2, \dots be a sequence of i.i.d. $\mathcal{N}(0, 1)$ random variables. Then for any integers $l \neq 0$ and $k > 0$, the following exponential probability bound holds for any $x > 0$:*

$$\mathbb{P} \left(\left| \sum_{t=1}^k Z_t Z_{t+l} \right| > kx \right) \leq 2 \exp \left(-\frac{1}{8} \frac{kx^2}{4+x} \right). \quad (\text{S4.10})$$

Proof. We will show that $\mathbb{P} \left(\sum_{t=1}^k Z_t Z_{t+l} > kx \right) \leq \exp \left(-\frac{1}{8} \frac{kx^2}{4+x} \right)$, which would then imply (S4.10) by symmetry. By Markov's inequality, for any $x > 0$ and $\lambda > 0$, it holds that

$$\mathbb{P} \left(\sum_{t=1}^k Z_t Z_{t+l} > kx \right) \leq \exp(-kx\lambda) \mathbb{E} \exp \left(\lambda \sum_{t=1}^k Z_t Z_{t+l} \right).$$

By the convexity of $y \mapsto \exp(\lambda y)$ for any $\lambda > 0$, Theorem 1 in Vershynin (2011) implies

$$\mathbb{E} \exp \left(\lambda \sum_{t=1}^k Z_t Z_{t+l} \right) \leq \mathbb{E} \exp \left(4\lambda \sum_{t=1}^k Z_t \tilde{Z}_t \right),$$

where $\tilde{Z}_1, \dots, \tilde{Z}_k$ are independent copies of Z_1, \dots, Z_k . Using the independence and by direct computation (see also Craig (1936)), we get

$$\mathbb{E} \exp \left(4\lambda \sum_{t=1}^k Z_t \tilde{Z}_t \right) = \left(\mathbb{E} \exp \left(4\lambda Z_1 \tilde{Z}_1 \right) \right)^k = \left(\mathbb{E} \exp \left(8\lambda^2 \tilde{Z}_1^2 \right) \right)^k = (1 - 16\lambda^2)^{-\frac{1}{2}k}$$

provided that $0 < \lambda < \frac{1}{4}$, therefore $\mathbb{P} \left(\sum_{t=1}^k Z_t Z_{t+l} > kx \right) \leq \exp \left(-kx\lambda - \frac{k}{2} \log(1 - 16\lambda^2) \right)$.

Taking $\lambda = \frac{-2 + \sqrt{4+x^2}}{4x}$ minimises the right-hand side of this inequality. With this value of

λ and using $\log(x) \leq x - 1$, we have

$$\begin{aligned} \mathbb{P} \left(\sum_{t=1}^k Z_t Z_{t+l} > kx \right) &\leq \exp \left(\frac{k}{4} \left(2 - \sqrt{x^2 + 4} + 2 \log \left(\frac{1}{4} (\sqrt{x^2 + 4} + 2) \right) \right) \right) \\ &\leq \exp \left(\frac{k}{4} \left(2 - \sqrt{x^2 + 4} + \frac{1}{2} (\sqrt{x^2 + 4} + 2) - 2 \right) \right) \\ &= \exp \left(\frac{k}{8} (2 - \sqrt{x^2 + 4}) \right) = \exp \left(-\frac{1}{8} \frac{kx^2}{2 + \sqrt{x^2 + 4}} \right) \\ &\leq \exp \left(-\frac{1}{8} \frac{kx^2}{4 + x} \right), \end{aligned}$$

which completes the proof. □

Lemma S4.5 (Lemma 1 in Laurent and Massart (2000)) *Let Z_1, Z_2, \dots be a sequence of i.i.d. $\mathcal{N}(0, 1)$ random variables. For any integer $k > 0$ and $x > 0$, the following*

exponential probability bounds hold

$$\mathbb{P} \left(\sum_{t=1}^k Z_t^2 \geq k + 2\sqrt{kx} + 2x \right) \leq \exp(-x), \quad (\text{S4.11})$$

$$\mathbb{P} \left(\sum_{t=1}^k Z_t^2 \leq k - 2\sqrt{kx} \right) \leq \exp(-x). \quad (\text{S4.12})$$

Proof of Theorem 2.1. For $\mathbf{C}_T = \sum_{t=1}^{T-1} \mathbf{Y}_t \mathbf{Y}_t'$ and $\mathbf{A}_T = \sum_{t=1}^{T-1} \varepsilon_{t+1} \mathbf{Y}_t$, we have $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = \mathbf{C}_T^{-1} \mathbf{A}_T$. Here the distribution of $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}$ is invariant to the value of σ . As such, in the following, we assume $\sigma = 1$ for notational convenience. Consequently,

$$\left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right\| \leq \lambda_{\max}(\mathbf{C}_T^{-1}) \|\mathbf{A}_T\| = \lambda_{\min}^{-1}(\mathbf{C}_T) \|\mathbf{A}_T\|, \quad (\text{S4.13})$$

where $\lambda_{\min}(\mathbf{M})$ and $\lambda_{\max}(\mathbf{M})$ denote, respectively, the smallest and the largest eigenvalues of a symmetric matrix \mathbf{M} . To provide an upper bound on $\left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right\|$ given in Theorem 2.1, we will bound $\lambda_{\min}(\mathbf{C}_T)$ from below and $\|\mathbf{A}_T\|$ from above, working on a set whose probability is large.

In the calculations below, we will repeatedly use the following representation of \mathbf{Y}_t , which follows from applying (S4.5) recursively:

$$\mathbf{Y}_t = \mathbf{B}^t \mathbf{Y}_0 + \sum_{j=1}^t \varepsilon_{t-j+1} \mathbf{B}^{j-1} \mathbf{u}, \quad t = 1, \dots, T. \quad (\text{S4.14})$$

In addition, to improve the presentational aspect of the proof, here we shall take $\mathbf{Y}_0 = \mathbf{0}$.

All the results would go through (with minor modifications to handle the extra terms) if

one instead assumes that \mathbf{Y}_0 is a realization from a stationary solution.

In the arguments below, we will show result more specific than (2.5), i.e.

$$\|\mathbf{A}_T\| \leq \left(32\bar{b}^{-2}\sqrt{1 + \|\boldsymbol{\beta}\|^2}\right) p \log(T) \sqrt{(1 + \log(T + p))T}, \quad (\text{S4.15})$$

$$\lambda_{\min}(\mathbf{C}_T) \geq \bar{b}^{-2} \left(T - p(1 + 32 \log(T)\sqrt{T})\right), \quad (\text{S4.16})$$

on the event

$$\mathcal{E}_T = \mathcal{E}_T^{(1)} \cap \mathcal{E}_T^{(2)} \cap \mathcal{E}_T^{(3)}, \quad (\text{S4.17})$$

where

$$\begin{aligned} \mathcal{E}_T^{(1)} &= \bigcap_{1 \leq i < j \leq p} \left\{ \left| \sum_{t=1}^{T - \max(i, j)} \varepsilon_t \varepsilon_{t+|i-j|} \right| < 32 \log(T) \sqrt{T - \max(i, j)} \right\}, \\ \mathcal{E}_T^{(2)} &= \bigcap_{j=1}^T \left\{ \left| \sum_{t=1}^{T-j} \varepsilon_t \varepsilon_{t+j} \right| < 32 \log(T) \sqrt{T - j} \right\}, \\ \mathcal{E}_T^{(3)} &= \left\{ \sum_{t=1}^{T-p} \varepsilon_t^2 > T - p - 2\sqrt{\log(T)(T - p)} \right\}. \end{aligned}$$

Finally, we will demonstrate that \mathcal{E}_T satisfies

$$\mathbb{P}(\mathcal{E}_T) \geq 1 - \frac{5}{T}. \quad (\text{S4.18})$$

Thus, (S4.13), (S4.15), (S4.16) and (S4.18) combined together imply the statement of

Theorem 2.1. The remaining part of the proof is split into three parts, in which we show (S4.15), (S4.16) and (S4.18) in turn.

Upper bound for $\|\mathbf{A}_T\|$. The Euclidean norm satisfies $\|\mathbf{A}_T\| = \sup_{\mathbf{v} \in \mathbb{R}^p, \|\mathbf{v}\|=1} |\langle \mathbf{v}, \mathbf{A}_T \rangle|$, therefore we consider inner products $\langle \mathbf{v}, \mathbf{A}_T \rangle$ where $\mathbf{v} \in \mathbb{R}^p$ is any unit vector. By (S4.14),

$$\begin{aligned} \langle \mathbf{v}, \mathbf{A}_T \rangle &= \sum_{t=1}^{T-1} \langle \mathbf{v}, \mathbf{Y}_t \rangle \varepsilon_{t+1} = \sum_{t=1}^{T-1} \sum_{j=1}^t \langle \mathbf{v}, \mathbf{B}^{j-1} \mathbf{u} \rangle \varepsilon_{t-j+1} \varepsilon_{t+1} \\ &= \sum_{j=1}^{T-1} \langle \mathbf{v}, \mathbf{B}^{j-1} \mathbf{u} \rangle a_j, \end{aligned}$$

where $a_j = \sum_{t=j}^{T-1} \varepsilon_{t-j+1} \varepsilon_{t+1} = \sum_{t=1}^{T-j} \varepsilon_t \varepsilon_{t+j}$.

Lemma S4.2 and Lemma S4.3 applied to the right-hand side of the above equation yield

$$\begin{aligned} \sum_{j=1}^{T-1} \langle \mathbf{v}, \mathbf{B}^{j-1} \mathbf{u} \rangle a_j &= \frac{1}{2\pi i} \int_{\mathbb{T}} \left(\sum_{j=1}^{T-1} z^{j-1} a_j \right) \langle \mathbf{v}, (z\mathbf{I}_p - \mathbf{B})^{-1} \mathbf{u} \rangle dz \\ &= \frac{1}{2\pi i} \int_{\mathbb{T}} \left(\sum_{j=1}^{T-1} z^{j-1} a_j \right) \left(\sum_{j=1}^p z^{p-j} v_j \right) q(z) dz \\ &= \frac{1}{2\pi i} \int_{\mathbb{T}} \left(\sum_{j=0}^{T+p-1} z^j c_j \right) q(z) dz, \end{aligned}$$

where $q(z) = (z^p b(z^{-1}))^{-1}$ and $c_j = \sum_{i=0}^j a_{i+1} v_{p-j+i}$. Integrating by parts, we get

$$\frac{1}{2\pi i} \int_{\mathbb{T}} \left(\sum_{j=0}^{T+p-1} z^j c_j \right) q(z) dz = -\frac{1}{2\pi i} \int_{\mathbb{T}} \left(\sum_{j=0}^{T+p-1} z^{j+1} \frac{c_j}{j+1} \right) q'(z) dz,$$

where $q'(\cdot)$ is the derivative of $q(\cdot)$. Combining the calculations above and using the fact

that $\mathbb{T} = \{z \in \mathbb{C} : |z| = 1\}$, Cauchy's inequality and Lemma S4.1, we obtain

$$\left| \sum_{j=1}^{T-1} \langle \mathbf{v}, \mathbf{B}^{j-1} \mathbf{u} \rangle a_j \right| \leq \sqrt{\sum_{j=0}^{T+p-1} \left(\frac{c_j}{j+1} \right)^2} \sqrt{\int_{\mathbb{T}} |q'(z)|^2 dm(z)}, \quad (\text{S4.19})$$

where we recall that $dm(z) = \frac{d|z|}{2\pi}$. To further bound the first term on the right-hand side of (S4.19), we recall that on the event \mathcal{E}_T coefficients $|a_j| \leq 32 \log(T) \sqrt{T}$, hence

$$\begin{aligned} \sqrt{\sum_{j=0}^{T+p-1} \left(\frac{c_j}{j+1} \right)^2} &= \sqrt{\sum_{j=0}^{T+p-1} \frac{1}{(j+1)^2} \left(\sum_{i=0}^j a_{i+1} v_{p-j+i} \right)^2} \\ &\leq \max_{j=0, \dots, T+p-1} |a_j| \sqrt{\sum_{j=0}^{T+p-1} \frac{1}{(j+1)^2} \left(\sum_{i=0}^j |v_{p-j+i}| \right)^2} \\ &\leq 32 \log(T) \sqrt{T} \sqrt{\sum_{j=0}^{T+p-1} \frac{j+1}{(j+1)^2}} \\ &\leq 32 \log(T) \sqrt{(1 + \log(T+p))T}. \end{aligned}$$

For the second term in (S4.19), we calculate the derivative

$$q'(z) = -\frac{pz^{p-1} - \sum_{j=1}^p (p-j)\beta_j z^{p-j-1}}{(z^p b(z^{-p}))^2}$$

and use Lemma S4.1 to bound

$$\begin{aligned}
\sqrt{\int_{\mathbb{T}} |q'(z)|^2 dm(z)} &= \sqrt{\int_{\mathbb{T}} \left| \frac{pz^{p-1} - \sum_{j=1}^p (p-j)\beta_j z^{p-j}}{(z^p b(z^{-p}))^2} \right|^2 dm(z)} \\
&\leq \frac{\sqrt{\int_{\mathbb{T}} \left| pz^{p-1} - \sum_{j=1}^p (p-j)\beta_j z^{p-j} \right|^2 dm(z)}}{\min_{|z|=1} |(z^p b(z^{-p}))|^2} \\
&= \underline{b}^{-2} \sqrt{\left(p^2 + \sum_{j=1}^p (p-j)^2 \beta_j^2 \right)} \leq \underline{b}^{-2} p \sqrt{1 + \|\boldsymbol{\beta}\|^2}.
\end{aligned}$$

Combining the bounds on the two terms, we obtain

$$\sum_{j=1}^{T-1} \langle \mathbf{v}, \mathbf{B}^{j-1} \mathbf{u} \rangle a_j \leq \left(32 \underline{b}^{-2} \sqrt{1 + \|\boldsymbol{\beta}\|^2} \right) p \log(T) \sqrt{(1 + \log(T+p))T}.$$

Taking supremum over $\mathbf{v} \in \mathbb{R}^p$ such that $\|\mathbf{v}\| = 1$ proves (S4.15).

Lower bound for $\lambda_{\min}(\mathbf{C}_T)$. Let $\mathbf{v} = (v_1, \dots, v_p)'$ be a unit vector in \mathbb{R}^p . We begin the proof by establishing the following inequality

$$\langle \mathbf{v}, \mathbf{C}_T \mathbf{v} \rangle \geq \bar{b}^{-2} \sum_{i,j=1}^p v_i v_j \sum_{t=1}^{T-1} \varepsilon_{t-j+1} \varepsilon_{t-i+1}, \tag{S4.20}$$

where $\varepsilon_t = 0$ for $t \leq 0$ and $\bar{b} = \max_{z \in \mathbb{T}} |b(z)|$. By Lemma S4.1 and (S4.14), we rewrite

the quadratic form on the left-hand side of (S4.20) to

$$\langle \mathbf{v}, \mathbf{C}_T \mathbf{v} \rangle = \sum_{t=1}^{T-1} \langle \mathbf{v}, \mathbf{Y}_t \rangle^2 \quad (\text{S4.21})$$

$$= \int_{\mathbb{T}} \left| \sum_{t=1}^{T-1} \left\langle \mathbf{v}, \sum_{j=1}^t \varepsilon_j \mathbf{B}^{t-j} \mathbf{u} \right\rangle z^t \right|^2 dm(z) \quad (\text{S4.22})$$

$$= \int_{\mathbb{T}} \left| \sum_{t=1}^{T-1} \sum_{j=1}^{T-1} \varepsilon_j \omega_{t-j} z^t \right|^2 dm(z) \quad (\text{S4.23})$$

where $\omega_j = \langle \mathbf{v}, \mathbf{B}^j \mathbf{u} \rangle$ for $j \geq 0$, $\omega_j = 0$ for $j < 0$. Changing the order of summation and by a simple substitution we get

$$\sum_{t=1}^{T-1} \sum_{j=1}^{T-1} \varepsilon_j \omega_{t-j} z^t = \sum_{j=1}^{T-1} \varepsilon_j z^j \sum_{t=1}^{T-1} \omega_{t-j} z^{t-j} = \sum_{j=1}^{T-1} \varepsilon_j z^j \sum_{t=0}^{T-j-1} \omega_t z^t. \quad (\text{S4.24})$$

Using the definition of ω_j , the fact that all eigenvalues of \mathbf{B} have modulus strictly lower than one and Lemma S4.3, (S4.24) simplifies to

$$\begin{aligned} \sum_{j=1}^{T-1} \varepsilon_j z^j \sum_{t=0}^{T-j-1} \omega_t z^t &= \sum_{j=1}^{T-1} \varepsilon_j z^j \langle \mathbf{v}, (\mathbf{I}_p - (\mathbf{B}z)^{T-j})(\mathbf{I}_p - \mathbf{B}z)^{-1} \mathbf{u} \rangle \\ &= \sum_{j=1}^{T-1} \varepsilon_j \left(z^j \langle \mathbf{v}, (\mathbf{I}_p - \mathbf{B}z)^{-1} \mathbf{u} \rangle - z^T \langle \mathbf{B}^{T-j} \mathbf{v}, (\mathbf{I}_p - \mathbf{B}z)^{-1} \mathbf{u} \rangle \right) \\ &= b(z)^{-1} \sum_{j=1}^{T-1} \varepsilon_j \left(z^j v(z) - z^T w_j(z) \right), \end{aligned}$$

where $v(z) = \sum_{k=1}^p v_k z_{k-1}$ and $w_j(z) = \sum_{k=1}^p (\mathbf{B}^{T-j} v)_k z^{k-1}$ for $j = 0, \dots, T-1$. The

equation above, (S4.21) and (S4.24) combined together imply the following inequality

$$\begin{aligned} \langle \mathbf{v}, \mathbf{C}_T \mathbf{v} \rangle &= \int_{\mathbb{T}} \left| b(z)^{-1} \sum_{j=1}^{T-1} \varepsilon_j (z^j v(z) - z^T w_j(z)) \right|^2 dm(z) \\ &\geq \bar{b}^{-2} \int_{\mathbb{T}} \left| \sum_{j=1}^{T-1} \varepsilon_j (z^j v(z) - z^T w_j(z)) \right|^2 dm(z). \end{aligned}$$

Observe that $\sum_{j=1}^{T-1} \varepsilon_j (z^j v(z) - z^T w_j(z)) = \sum_{j=1}^{T-1} \varepsilon_j (z^j v(z) - z^T w_j(z)) = \sum_{t=1}^{T+p-1} c_t z^t$ is a trigonometric polynomial, therefore by Lemma S4.1 and simple algebra

$$\begin{aligned} \int_{\mathbb{T}} \left| \sum_{j=1}^{T-1} \varepsilon_j (z^j v(z) - z^T w_j(z)) \right|^2 dm(z) &= \sum_{t=1}^{T+p-1} |c_t|^2 \geq \sum_{t=1}^{T-1} |c_t|^2 = \sum_{t=1}^{T-1} \left(\sum_{j=1}^p v_j \varepsilon_{t-j+1} \right)^2 = \\ &= \sum_{i,j=1}^p v_j v_i \sum_{t=1}^{T-1} \varepsilon_{t-j+1} \varepsilon_{t-i+1}, \end{aligned}$$

which proves (S4.20).

We are now in a position to bound $\langle \mathbf{v}, \mathbf{C}_T \mathbf{v} \rangle$ from below. Rearranging terms in (S4.20) yields

$$\begin{aligned} \langle \mathbf{v}, \mathbf{C}_T \mathbf{v} \rangle &\geq \bar{b}^{-2} \left(\sum_{i=1}^p v_i^2 \sum_{t=1}^{n-i} \varepsilon_t^2 + \sum_{1 \leq i < j \leq p} v_i v_j \sum_{t=1}^{T-\max(i,j)} \varepsilon_t \varepsilon_{t+|j-i|} \right) \\ &\geq \bar{b}^{-2} \left(\sum_{t=1}^{T-p} \varepsilon_t^2 \sum_{i=1}^p v_i^2 - \max_{1 \leq i < j \leq p} \left| \sum_{t=1}^{T-\max(i,j)} \varepsilon_t \varepsilon_{t+|j-i|} \right| \left(\left(\sum_{i=1}^p |v_i| \right)^2 - \sum_{i=1}^p v_i^2 \right) \right) \\ &\geq \bar{b}^{-2} \left(\sum_{t=1}^{T-p} \varepsilon_t^2 - (p-1) \max_{1 \leq i < j \leq p} \left| \sum_{t=1}^{T-\max(i,j)} \varepsilon_t \varepsilon_{t+|j-i|} \right| \right). \end{aligned}$$

Recalling the definition of \mathcal{E}_T , we conclude that on this event

$$\begin{aligned} \langle \mathbf{v}, \mathbf{C}_T \mathbf{v} \rangle &\geq \bar{b}^{-2} \left(T - p - 2\sqrt{\log(T)(T-p)} - (p-1)32\log(T)\sqrt{T} \right) \\ &\geq \bar{b}^{-2} \left(T - p(1 + 32\log(T)\sqrt{T}) \right). \end{aligned}$$

Taking infimum over $\mathbf{v} \in \mathbb{R}^p$ such that $\|\mathbf{v}\| = 1$ in the inequality above proves (S4.16).

Lower bound for $\mathbb{P}(\mathcal{E}_T)$. Recalling (S4.17) and using a simple Bonferroni bound, we get

$$\begin{aligned} \mathbb{P}(\mathcal{E}_T^c) &\leq p^2 \max_{1 \leq i < j \leq p} \mathbb{P} \left(\left| \sum_{t=1}^{T-\max(i,j)} \varepsilon_t \varepsilon_{t+|i-j|} \right| \geq 32\log(T)\sqrt{T-\max(i,j)} \right) \\ &\quad + T \max_{1 \leq j \leq T} \mathbb{P} \left(\left| \sum_{t=1}^{T-j} \varepsilon_t \varepsilon_{t+j} \right| < 32\log(T)\sqrt{T-j} \right) \\ &\quad + \mathbb{P} \left(\sum_{t=1}^{T-p} \varepsilon_t^2 > T - p - 2\sqrt{\log(T)(T-p)} \right) \\ &:= p^2 \max_{1 \leq i < j \leq p} P_{i,j}^{(1)} + T \max_{1 \leq j \leq T} P_j^{(2)} + P^{(3)}. \end{aligned}$$

Lemma S4.4 implies that

$$\begin{aligned} P_{i,j}^{(1)} &\leq 2 \exp \left(-\frac{1}{8} \frac{(32\log(T))^2}{4 + (\sqrt{T-\max(i,j)})^{-1}32\log(T)} \right) \leq 2 \exp(-2\log(T)) = \frac{2}{T^2}, \\ P_j^{(2)} &\leq 2 \exp \left(-\frac{1}{8} \frac{(32\log(T))^2}{4 + (\sqrt{T-j})^{-1}32\log(T)} \right) \leq 2 \exp(-2\log(T)) = \frac{2}{T^2}. \end{aligned}$$

Moreover, by Lemma S4.5, $P^{(3)} \leq \exp(-\log(T)) = \frac{1}{T}$, hence, given that $p^2 < T$, we have

$\mathbb{P}(\mathcal{E}_T^c) \leq \frac{5}{T}$, which completes the proof. \square

S4.3 Proof of Theorem 2.2

In the proof below, we shall focus on the case where F_T^M consists of randomly drawn intervals (which is what Algorithm 2 does when p is large). For the case where all sub-intervals of $[1, p]$ are used, the same arguments would go through, because Algorithm 2 then produces a larger set F_T^M compared to the approach of random drawing.

We now split the proof into four steps.

Step 1. Consider the event $\left\{ \left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right\| \leq \kappa_1 (\underline{b}/\bar{b})^2 \|\boldsymbol{\beta}\| \frac{p \log(T) \sqrt{\log(T+p)}}{\sqrt{T - \kappa_2 p \log(T)}} \right\}$ where κ_1, κ_2 are as in Theorem 2.1. Assumption (A3) implies that \underline{b}/\bar{b} and $\|\boldsymbol{\beta}\|$ are bounded from above by constants. Furthermore, by Assumption (A2), $p \leq c_1 T^\theta$, which implies that

$$\kappa_1 (\underline{b}/\bar{b})^2 \|\boldsymbol{\beta}\| \frac{p \log(T) \sqrt{\log(T+p)}}{\sqrt{T - \kappa_2 p \log(T)}} \leq c_3 T^{\theta-1/2} (\log(T))^{3/2} = c_3 \lambda_T =: \lambda_T \quad (\text{S4.25})$$

for some constant $c_3 > 0$ and a sufficiently large T . Define now

$$A_T = \left\{ \left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right\| \leq \lambda_T \right\} \quad (\text{S4.26})$$

By Theorem 2.1,

$$\mathbb{P}(A_T) \geq \mathbb{P} \left(\left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \right\| \leq \kappa_1 (\underline{b}/\bar{b})^2 \|\boldsymbol{\beta}\| \frac{p \log(T) \sqrt{\log(T+p)}}{\sqrt{T - \kappa_2 p \log(T)}} \right) \geq 1 - \kappa_3 T^{-1}, \quad (\text{S4.27})$$

for some constant $\kappa_3 > 0$.

Step 2. For $j = 1, \dots, q$, define the intervals

$$\mathcal{I}_j^L = (\tau_j - \delta_T/3, \tau_j - \delta_T/6) \quad (\text{S4.28})$$

$$\mathcal{I}_j^R = (\tau_j + \delta_T/6, \tau_j + \delta_T/3) \quad (\text{S4.29})$$

Recall that F_T^M is the set of M randomly drawn intervals with endpoints in $\{1, \dots, p\}$.

Denote by $[s_1, e_1], \dots, [s_M, e_M]$ the elements of F_T^M and let

$$D_T^M = \left\{ \forall j = 1, \dots, q, \exists k \in \{1, \dots, M\}, \text{ s.t. } s_k \times e_k \in \mathcal{I}_j^L \times \mathcal{I}_j^R \right\}. \quad (\text{S4.30})$$

We have that

$$\begin{aligned} \mathbb{P}((D_T^M)^c) &\leq \sum_{j=1}^q \prod_{m=1}^M \left(1 - \mathbb{P}(s_m \times e_m \in \mathcal{I}_j^L \times \mathcal{I}_j^R) \right) \\ &\leq q \left(1 - \frac{\delta_T^2}{6^2 p^2} \right)^M \leq \frac{p}{\delta_T} \left(1 - \frac{\delta_T^2}{36 p^2} \right)^M. \end{aligned}$$

Therefore, $\mathbb{P}(A_T \cap D_T^M) \geq 1 - \kappa_3 T^{-1} - p \delta_T^{-1} (1 - \delta_T^2 p^{-2}/36)^M \rightarrow 1$. Note that the same conclusion still holds if F_T^M contains all the intervals with endpoints in $\{1, \dots, p\}$. In the remainder of the proof, assume that A_T and D_T^M all hold.

Note that Assumption (A4) implies that there exists $\underline{c} > 0$ such that $\delta_T^{1/2} \underline{\alpha}_T > \underline{c} \lambda_T$ for all sufficiently large T . We are now in the position to specify the constants explicitly

as

$$C_1 = 2\sqrt{C_3} + c_3, \quad C_2 = \frac{1}{\sqrt{6}} - \frac{1}{\underline{c}}, \quad C_3 = (4\sqrt{2} + 6)c_3^2,$$

where c_3 is in Equation (S4.25).

Step 3. We focus on a generic interval $[s, e]$ such that

$$\exists j \in \{1, \dots, q\}, \exists k \in \{1, \dots, M\}, \text{ s.t. } [s_k, e_k] \subset [s, e] \text{ and } s_k \times e_k \in \mathcal{I}_j^L \times \mathcal{I}_j^R. \quad (\text{S4.31})$$

Fix such an interval $[s, e]$ and let $j \in \{1, \dots, q\}$ and $k \in \{1, \dots, M\}$ be such that (S4.31) is satisfied. Let $b_k^* = \operatorname{argmax}_{s_k \leq b \leq e_k} \mathcal{C}_{s_k, e_k}^b(\hat{\beta})$. By construction, $[s_k, e_k]$ satisfies $\tau_j - s_k + 1 \geq \delta_T/6$ and $e_k - \tau_j > \delta_T/6$. Let

$$\begin{aligned} \mathcal{M}_{s,e} &= \{m : [s_m, e_m] \in F_T^M, [s_m, e_m] \subset [s, e]\}, \\ \mathcal{O}_{s,e} &= \{m \in \mathcal{M}_{s,e} : \max_{s_m \leq b < e_m} \mathcal{C}_{s_m, e_m}^b(\hat{\beta}) > \zeta_T\}. \end{aligned}$$

Our first aim is to show that $\mathcal{O}_{s,e}$ is non-empty. This follows from Lemma 2 in Baranowski et al. (2019), the Cauchy–Schwarz inequality, and the calculation below, as

$$\begin{aligned} \mathcal{C}_{s_k, e_k}^{b_k^*}(\hat{\beta}) &\geq \mathcal{C}_{s_k, e_k}^{\tau_j}(\hat{\beta}) \\ &\geq \mathcal{C}_{s_k, e_k}^{b_k^*}(\beta) - \lambda_T \geq \left(\frac{\delta_T}{6}\right)^{1/2} |\alpha_j \tau_j^{-1}| - \lambda_T \geq \left(\frac{\delta_T}{6}\right)^{1/2} \underline{\alpha}_T - \lambda_T \\ &= \left(\frac{1}{\sqrt{6}} - \frac{\lambda_T}{\delta_T^{1/2} \underline{\alpha}_T}\right) \delta_T^{1/2} \underline{\alpha}_T \geq \left(\frac{1}{\sqrt{6}} - \frac{1}{\underline{c}}\right) \delta_T^{1/2} \underline{\alpha}_T = C_2 \delta_T^{1/2} \underline{\alpha}_T > \zeta_T. \end{aligned}$$

Let $m^* = \operatorname{argmin}_{m \in \mathcal{O}_{s,e}} (e_m - s_m + 1)$ and $b^* = \operatorname{argmax}_{s_{m^*} \leq b < e_{m^*}} \mathcal{C}_{s_{m^*}, e_{m^*}}^b(\hat{\beta})$. Observe that $[s_{m^*}, e_{m^*})$ must contain at least one change in $\hat{\beta}$. Indeed, if this were not the case, we would have $\mathcal{C}_{s_{m^*}, e_{m^*}}^b(\beta) = 0$ and

$$\mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\hat{\beta}) = |\mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\hat{\beta}) - \mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\beta)| \leq \lambda_T < \frac{C_1}{c_3} \lambda_T = C_1 \underline{\lambda}_T \leq \zeta_T,$$

which contradicted $\mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\hat{\beta}) > \zeta_T$. On the other hand, $[s_{m^*}, e_{m^*})$ cannot contain more than one change-points, because $e_{m^*} - s_{m^*} + 1 \leq e_k - s_k + 1 \leq \delta_T$.

Without loss of generality, assume $\tau_j \in [s_{m^*}, e_{m^*})$. Let $\eta_L = \tau_j - s_{m^*} + 1$, $\eta_R = e_{m^*} - \tau_j$ and $\eta_T = (C_1/c_3 - 1)^2 \alpha_j^2 \tau_j^{-2} \lambda_T^2$. We claim that $\min(\eta_L, \eta_R) > \eta_T$, because otherwise $\min(\eta_L, \eta_R) \leq \eta_T$ and Lemma 2 in Baranowski et al. (2019) would have implied

$$\begin{aligned} \mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\hat{\beta}) &\leq \mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\beta) + \lambda_T \leq \mathcal{C}_{s_{m^*}, e_{m^*}}^{\tau_j}(\beta) + \lambda_T \leq \eta_T^{1/2} |\alpha_j \tau_j^{-1}| + \lambda_T \\ &= (C_1/c_3 - 1 + 1) \lambda_T = C_1 \underline{\lambda}_T < \zeta_T, \end{aligned}$$

which contradicted $\mathcal{C}_{s_{m^*}, e_{m^*}}^{b^*}(\hat{\beta}) > \zeta_T$.

We are now in the position to prove $|b^* - \tau_j| \leq C_3 \underline{\lambda}_T \alpha_T^{-2}$. Our aim is to find ϵ_T such that for any $b \in \{s_{m^*}, s_{m^*} + 1, \dots, e_{m^*} - 1\}$ with $|b - \tau_j| > \epsilon_T$, we always have

$$\left\{ \mathcal{C}_{s_{m^*}, e_{m^*}}^{\tau_j}(\hat{\beta}) \right\}^2 - \left\{ \mathcal{C}_{s_{m^*}, e_{m^*}}^b(\hat{\beta}) \right\}^2 > 0. \quad (\text{S4.32})$$

This would then imply that $|b^* - \tau_j| \leq \epsilon_T$. By expansion and rearranging the terms, we

see that (S4.32) is equivalent to

$$\begin{aligned} & \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle^2 - \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle^2 > \langle \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle^2 - \langle \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle^2 \\ & + 2 \left\langle \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle - \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle \right\rangle. \end{aligned} \quad (\text{S4.33})$$

Here $\boldsymbol{\psi}_{s,e}^b$ (with $1 \leq s < b < e \leq p$) is a p -dimensional vector, with its s -th to b -th component being $\sqrt{\frac{e-b}{(e-s+1)(b-s+1)}}$, its $b+1$ -th to e -th component being $\sqrt{\frac{b-s+1}{(e-s+1)(e-b)}}$, and the remaining elements being 0. In the following, we assume that $b \geq \tau_j$. The case that $b < \tau_j$ can be handled in a similar fashion. By Lemma 4 in Baranowski et al. (2019), we have

$$\begin{aligned} \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle^2 - \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle^2 &= (\mathcal{C}_{s^*, e^*}^{\tau_j}(\boldsymbol{\beta}))^2 - (\mathcal{C}_{s_m^*, e_m^*}^b(\boldsymbol{\beta}))^2 \\ &= \frac{|b - \tau_j| \eta_L}{|b - \tau_j| + \eta_L} (\alpha_j \tau_j^{-1})^2 =: \kappa. \end{aligned}$$

In addition, since we assume event A_T ,

$$\begin{aligned} & \langle \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle^2 - \langle \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle^2 \leq \lambda_T^2, \\ & 2 \left\langle \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle - \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle \right\rangle \\ & \leq 2 \|\boldsymbol{\psi}_{s_m^*, e_m^*}^b \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^b \rangle - \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \langle \boldsymbol{\beta}, \boldsymbol{\psi}_{s_m^*, e_m^*}^{\tau_j} \rangle\|_2 \lambda_T = 2\kappa^{1/2} \lambda_T, \end{aligned}$$

where the final equality is also implied by Lemma 4 in Baranowski et al. (2019). Consequently, (S4.33) can be deduced from the stronger inequality $\kappa - 2\lambda_T \kappa^{1/2} - \lambda_T^2 > 0$.

This quadratic inequality is implied by $\kappa > (\sqrt{2} + 1)^2 \lambda_T^2$, and could be restricted further to

$$\frac{2|b - \tau_j| \eta_L}{|b - \tau_j| + \eta_L} \geq \min(|b - \tau_j|, \eta_L) > (4\sqrt{2} + 6)(\alpha_j \tau_j^{-1})^{-2} \lambda_T^2 = C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2. \quad (\text{S4.34})$$

But since

$$\eta_L \geq \eta_T = (C_1/c_3 - 1)^2 (\alpha_j \tau_j^{-1})^{-2} \lambda_T^2 = (2\sqrt{C_3}/c_3)^2 (\alpha_j \tau_j^{-1})^{-2} \lambda_T^2 > C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2,$$

we see that (S4.34) is implied by $|b - \tau_j| > C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2$. To sum up, $|b^* - \tau_j| > C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2$ would result in (S4.32), a contradiction. So we have proved that $|b^* - \tau_j| \leq C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2$.

Step 4. With the arguments above valid on the event $A_T \cap B_T \cap D_T^M$, we can now proceed with the proof of the theorem. At the start of Algorithm 1, we have $s = 1$ and $e = p$ and, provided that $q \geq 1$, condition (S4.31) is satisfied. Therefore the algorithm detects a change-point b^* in that interval such that $|b^* - \tau_j| \leq C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2$. By construction, we also have that $|b^* - \tau_j| < 2/3\delta_T$. This in turn implies that for all $l = 1, \dots, q$ such that $\tau_l \in [s, e]$ and $l \neq j$ we have either $\mathcal{I}_l^L, \mathcal{I}_l^R \subset [s, b^*]$ or $\mathcal{I}_l^L, \mathcal{I}_l^R \subset [b^* + 1, e]$. Therefore (S4.31) is satisfied within each segment containing at least one change-point. Note that before all q change points are detected, each change point will not be detected twice. To see this, we suppose that τ_j has already been detected by b , then for all intervals $[s_k, e_k] \subset [\tau_j - C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2 + 1, \tau_j - C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2 + 2/3\delta_T + 1] \cup$

$[\tau_j + C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2 - 2/3\delta_T, \tau_j + C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2]$, Lemma 2 in Baranowski et al. (2019), together with the definition of A_T , guarantee that

$$\begin{aligned} \max_{s_k \leq b < e} \mathcal{C}_{s_k, e_k}^b(\hat{\boldsymbol{\beta}}) &\leq \max_{s \leq b < e} \mathcal{C}_{s_k, e_k}^b(\boldsymbol{\beta}) + \lambda_T \\ &\leq \sqrt{C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2 \alpha_j \tau_j^{-1}} + \sqrt{C_3(\alpha_{j+1} \tau_{j+1}^{-1})^{-2} \underline{\lambda}_T^2 \alpha_{j+1} \tau_{j+1}^{-1}} + \lambda_T \\ &< (2\sqrt{C_3/c_3} + 1)\lambda_T = C_1 \underline{\lambda}_T < \zeta_T. \end{aligned}$$

Once all the change-points have been detected, we then only need to consider $[s_k, e_k]$ such that

$$[s_k, e_k] \subset [\tau_j - C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2 + 1, \tau_{j+1} + C_3(\alpha_{j+1} \tau_{j+1}^{-1})^{-2} \underline{\lambda}_T^2]$$

for $j = 1, \dots, q$. For such intervals, we have, by Lemmas 2 and 3 of Baranowski et al. (2019)

$$\begin{aligned} \max_{s_k \leq b < e_k} \mathcal{C}_{s_k, e_k}^b(\hat{\boldsymbol{\beta}}) &\leq \max_{s \leq b < e} \mathcal{C}_{s_k, e_k}^b(\boldsymbol{\beta}) + \lambda_T \\ &\leq \sqrt{C_3(\alpha_j \tau_j^{-1})^{-2} \underline{\lambda}_T^2 \alpha_j \tau_j^{-1}} + \sqrt{C_3(\alpha_{j+1} \tau_{j+1}^{-1})^{-2} \underline{\lambda}_T^2 \alpha_{j+1} \tau_{j+1}^{-1}} + \lambda_T \leq C_1 \underline{\lambda}_T < \zeta_T. \end{aligned}$$

Hence no further scales is detected and the algorithm terminates. \square

S4.4 Proof of Theorem 2.3

The proof of Theorem 2.3 is similar to that of Theorem 2.2. In the following, we shall still divide the proof into four steps as before, but focus on the main differences.

Step 1. Let $\{\rho_h : h \in \mathbb{Z}\}$ the true auto-correlation function of $\{X_t\}$ and $\hat{\rho}_h$ be its sample version (without de-meaning). Let $\boldsymbol{\rho} = (\rho_0, \dots, \rho_p)'$. First, we note that for $\alpha > 2$, the distribution of the innovations has finite second moment. It then follows from Anderson and Walker (1964) that $\hat{\boldsymbol{\rho}} - \boldsymbol{\rho} = O_p(T^{-1/2})$. The least-square estimator for AR(p) can be written as

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \sum_{i=p}^{T-1} X_i^2 & \cdots & \sum_{i=p}^{T-1} X_i X_{i-p+1} \\ & \ddots & \\ \sum_{i=1}^{T-p} X_i X_{i+p-1} & \cdots & \sum_{i=1}^{T-p} X_i^2 \end{bmatrix}_{p \times p}^{-1} \begin{bmatrix} \sum_{i=p}^{T-1} X_i X_{i+1} \\ \vdots \\ \sum_{i=1}^{T-p} X_i X_{i+p} \end{bmatrix}.$$

This is asymptotically equivalent to

$$\begin{bmatrix} \hat{\rho}_0 & \cdots & \hat{\rho}_{p-1} \\ & \ddots & \\ \hat{\rho}_{p-1} & \cdots & \hat{\rho}_0 \end{bmatrix}^{-1} \begin{bmatrix} \hat{\rho}_p \\ \vdots \\ \hat{\rho}_1 \end{bmatrix},$$

which converges to $\boldsymbol{\beta}$ at $O_p(T^{-1/2})$ in view of the Yule–Walker equations. Now for $0 < \alpha \leq 2$, despite infinite second moment in the innovations thus the time series, the auto-correlation function is still well-defined, in the sense of Davis and Resnick (1986). It follows from Hannan and Kanter (1977) that for any sufficiently small $\epsilon > 0$, $T^{1/\alpha-\epsilon} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\| \rightarrow 0$ in probability. See also Yohai and Maronna (1977) and Davis and Resnick (1986). In conclusion, we have that $T^{\max(1/2, 1/\alpha)-\epsilon} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\| \rightarrow 0$ in probability.

Steps 2 and 3. The following arguments are simpler, due to the fact that p is fixed. Because we go through all the intervals $[s, e]$ over $\{1, \dots, p\}$, we could see that under the event that $T^{\max(1/2, 1/\alpha) - \epsilon} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_\infty < 1$ (N.B. here the norm does not matter, as p is fixed), for any $j = 1, \dots, q$, and taking $C_1 = \sqrt{p}$ and $C_2 = 1/2$,

$$\max_{\tau_j \leq b < \tau_{j+1}} C_{\tau_j, \tau_{j+1}}^b(\boldsymbol{\beta}) \geq \underline{\alpha}_T / \sqrt{2} - 2 \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_\infty > C_2 \underline{\alpha}_T.$$

On the other hand, for all the intervals $[s, e]$ that do not include any of the change-points $\{\tau_1, \dots, \tau_q\}$,

$$\max_{s \leq b < e} C_{s, e}^b(\boldsymbol{\beta}) \leq \sqrt{p} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_\infty < C_1 T^{-\max(1/2, 1/\alpha) + \epsilon}.$$

Step 4. We shall now proceed with the proof under the event that $T^{\max(1/2, 1/\alpha) - \epsilon} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_\infty < 1$, which happens with probability one as $T \rightarrow \infty$. At the start of Algorithm 1, we have $s = 1$ and $e = p$. Since we pick the threshold $\zeta_T < C_2 \underline{\alpha}_T$, and we consider only the narrowest intervals with the corresponding contrasts (i.e. CUSUM-type statistic) over the threshold, we would end up considering all $[\tau_j, \tau_j + 1]$ for $j \in \{1, \dots, q\}$. Notice that before all the q change-points are detected, we would not consider other longer intervals, because of the nature of Algorithm 1. In addition, we will not consider intervals without any change because their corresponding contrast values would be below the threshold, as proved in the previous step. Once all the changes are detected, we then only need to consider the intervals located in between consecutive change-points, which all have corresponding contrast values smaller than $C_1 T^{\epsilon - \max(1/2, 1/\alpha)}$, thus the threshold ζ_T . Hence the algorithm would terminate with no further scales detected.

References

- A. Anastasiou, Y. Chen, H. Cho, and P. Fryzlewicz. **breakfast**: Methods for fast multiple change-point detection and estimation, 2021. URL <https://CRAN.R-project.org/package=breakfast>. R package version 2.2.
- T. W. Anderson and A. M. Walker. On the asymptotic distribution of the autocorrelations of a sample from a linear stochastic process. *The Annals of Mathematical Statistics*, 35:1296 – 1303, 1964.
- R. Baranowski, Y. Chen, and P. Fryzlewicz. Narrowest-over-threshold detection of multiple change points and change-point-like features. *Journal of the Royal Statistical Society Series B*, 81:649–672, 2019.
- H. Cho and P. Fryzlewicz. Multiple change point detection under serial dependence: Wild contrast maximisation and gappy Schwarz algorithm. *Preprint*, 2021.
- C. C. Craig. On the frequency function of xy . *The Annals of Mathematical Statistics*, 7: 1–15, 1936.
- R. Davis and S. Resnick. Limit theory for the sample covariance and correlation functions of moving averages. *The Annals of Statistics*, 14:533 – 558, 1986.
- J. Duoandikoetxea. *Fourier Analysis*, volume 29 of *Graduate Studies in Mathematics*. American Mathematical Society, 2001.

-
- P. Fearnhead and P. Clifford. On-line inference for hidden Markov models via particle filters. *Journal of the Royal Statistical Society Series B*, 65:887–899, 2003.
- P. Fryzlewicz. Wild binary segmentation for multiple change-point detection. *Annals of Statistics*, 42:2243–2281, 2014.
- E. J. Hannan and M. Kanter. Autoregressive processes with infinite variance. *Journal of Applied Probability*, 14:411–415, 1977.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28:1302–1338, 2000.
- J. O Ruanaidh and W. Fitzgerald. *Numerical Bayesian Methods Applied to Signal Processing*. Springer, 1996.
- R. Vershynin. A simple decoupling inequality in probability theory. Technical report, University of Michigan, 2011. URL <https://www.math.uci.edu/~rvershyn/papers/decoupling-simple.pdf>.
- V. J. Yohai and R. A. Maronna. Asymptotic behavior of least-squares estimates for autoregressive processes with infinite variances. *The Annals of Statistics*, 5:554–560, 1977.