

# Supplementary Material for Identifiable Autoregressive Variational Autoencoders for Nonlinear and Nonstationary Spatio-Temporal Blind Source Separation <sup>\*</sup>

Mika Sipilä<sup>1</sup> (✉), Klaus Nordhausen<sup>2</sup>, and Sara Taskinen<sup>1</sup>

<sup>1</sup> Department of Mathematics and Statistics, University of Jyväskylä, Finland

<sup>2</sup> Department of Mathematics and Statistics, University of Helsinki, Finland

## A Lemmas for the autoregressive exponential families

In this section, some useful Lemmas are given for univariate autoregressive exponential family distributions.

**Definition 3 (Autoregressive models).** *A generative model of  $x$  is considered to be autoregressive, if it can be written as*

$$x(\boldsymbol{\theta}^t) = \mu(\boldsymbol{\theta}^t) + \sum_{r=1}^R \gamma_r(\boldsymbol{\theta}^t) \left( x^{t-r} - \mu(\boldsymbol{\theta}^{t-r}) \right) + \omega(\boldsymbol{\theta}^t), \quad (14)$$

where  $\boldsymbol{\theta}^t \in \mathbb{R}^m$  is the parameter vector at time step  $t$ ,  $\mu$  is a trend function,  $\gamma_1, \dots, \gamma_R$  are the functions for autoregressive coefficients and  $\omega$  is white noise so that  $E(\omega(\boldsymbol{\theta}^t)) = 0$  and  $\text{Var}(\omega(\boldsymbol{\theta}^t)) < \infty$  for all  $t = 1, \dots, T$ , and  $\text{Cov}(\omega(\boldsymbol{\theta}^t), \omega(\boldsymbol{\theta}^{t'})) = 0$  for all  $t \neq t'$ . To ensure local weak-sense stationarity for each  $t$ , the (complex) roots  $y_i$  of the polynomial  $1 - \sum_{i=1}^R \gamma_i(\boldsymbol{\theta}_{t-i}) y^i$  must satisfy  $|y_i| > 1$ .

**Definition 4 (Autoregressive exponential family).** *Assume an autoregressive model defined by Definition 3. The univariate distribution  $p(x^t | \{x^{t-1:t-R}; \boldsymbol{\theta}^t\})$  belongs in univariate autoregressive exponential family, if its probability distribution can be written as*

$$p(x^t | \{x^{t-1:t-R}; \boldsymbol{\theta}^t\}) = \frac{Q(x^t, \{x^{t-1:t-R}\})}{Z(\boldsymbol{\theta}^t)} e^{\sum_{j=1}^k T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t)}, \quad (15)$$

where  $Q$  is a base measure,  $Z$  is a normalizing constant,  $T_1, \dots, T_k$  are sufficient statistics and  $\boldsymbol{\theta}^t$  is a parameter vector at time point  $t$ . The dimension  $k \in \{1, 2, \dots\}$  is assumed to be minimal, meaning that the distribution  $p$  cannot be written in form (15) using a smaller  $k' < k$ .

**Lemma 1.** *Consider autoregressive exponential family distribution. The components of sufficient statistics  $\mathbf{T}$  of the distribution are linearly independent. In other words, if there exists  $\boldsymbol{\alpha} \in \mathbb{R}^k$  so that  $\alpha_1 T_1(x^t, \{x^{t-1:t-R}\}) + \dots + \alpha_k T_k(x^t, \{x^{t-1:t-R}\}) = 0$ , then  $\boldsymbol{\alpha} = \mathbf{0}$ .*

*Proof:* Assume that the components of  $\mathbf{T}$  are not linearly independent. Then, there exists  $\boldsymbol{\alpha} \in \mathbb{R}^k$ ,  $\boldsymbol{\alpha} \neq \mathbf{0}$ , meaning that for some  $i \in \{1, \dots, k\}$ ,  $\alpha_i \neq 0$ . By reordering the indices, we can assume that  $\alpha_k \neq 0$ . Then, we can write  $T_k(x^t, \{x^{t-1:t-R}\}) = \sum_{j=1}^{k-1} \frac{\alpha_j}{\alpha_k} T_j(x^t, \{x^{t-1:t-R}\})$ . Let  $\lambda_j^*(\boldsymbol{\theta}^t) := (\lambda_j(\boldsymbol{\theta}^t) + \frac{\alpha_j}{\alpha_k} \lambda_k(\boldsymbol{\theta}^t))$ . Then, the term in the exponent of (15) can be written as

$$\sum_{j=1}^k T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) = \sum_{j=1}^{k-1} T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) + \sum_{j=1}^{k-1} \frac{\alpha_j}{\alpha_k} T_k(x^t, \{x^{t-1:t-R}\}) \quad (16)$$

$$= \sum_{j=1}^{k-1} T_j(x^t, \{x^{t-1:t-R}\}) \left( \lambda_j(\boldsymbol{\theta}^t) + \frac{\alpha_j}{\alpha_k} \lambda_k(\boldsymbol{\theta}^t) \right) \quad (17)$$

$$= \sum_{j=1}^{k-1} T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j^*(\boldsymbol{\theta}^t), \quad (18)$$

which contradicts the minimality of  $k$  in Definition 4. □

<sup>\*</sup> This research was supported by the Research Council of Finland (363261, 453691) and the Vilho, Yrjö and Kalle Väisälä foundation.

**Definition 5 (Strongly exponential autoregressive distributions).** *Exponential autoregressive distribution is considered strongly exponential if the following holds:*

$$(\exists \boldsymbol{\theta}^t \in \mathbb{R}^m \mid \forall x^t, \dots, x^{t-R} \in \mathcal{X}, \sum_{j=1}^k T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) = \text{const}) \implies l(\mathcal{X}) = 0 \text{ or } \boldsymbol{\lambda}(\boldsymbol{\theta}^t) = \mathbf{0}, \quad (19)$$

where  $l$  is a Lebesgue measure.

Definition 5 says that strongly exponential distribution has the exponential component in its expression almost surely, and the distribution can be reduced only to base measure and normalizing constant on a set of measure zero.

**Lemma 2.** *Consider a strongly exponential autoregressive family distribution whose sufficient statistics  $\mathbf{T}$  are differentiable almost everywhere. Then,  $T'_j \neq 0$  for all  $j = 1, \dots, k$  almost everywhere on  $\mathbb{R}$ .*

*Proof:* Assume that  $p$  is strongly exponential autoregressive distribution. Let  $\mathcal{X} = \cup_j \{x \in \mathbb{R}, T'_j(x) = 0\}$  and select any  $\boldsymbol{\theta}$  for which  $\boldsymbol{\lambda}(\boldsymbol{\theta}^t) \neq \mathbf{0}$ . Then, for all  $x \in \mathcal{X}$ , it holds that

$$\sum_{j=1}^k T'_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) = 0 \quad (20)$$

$$\implies \sum_{j=1}^k T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) = \text{const.} \quad (21)$$

By Definition 5, this means that  $l(\mathcal{X}) = 0$ . □

**Lemma 3.** *Consider a strongly exponential autoregressive family distribution of size  $k \geq 2$  so that the sufficient statistics  $\mathbf{T}$  are differentiable almost everywhere. Then, there exist  $k$  distinct points  $(x_1^t, \dots, x_1^{t-R}), \dots, (x_k^t, \dots, x_k^{t-R})$  such that the vectors  $\mathbf{T}'(x_1^t, \{x_1^{t-1:t-R}\}), \dots, \mathbf{T}'(x_k^t, \{x_k^{t-1:t-R}\})$  are linearly independent in  $\mathbb{R}^k$ .*

*Proof:* Suppose that for any choice of such  $k$  points, the vectors  $\mathbf{T}'(x_1^t, \{x_1^{t-1:t-R}\}), \dots, \mathbf{T}'(x_k^t, \{x_k^{t-1:t-R}\})$  are not linearly independent, meaning that there are a subspace of  $\mathbb{R}^k$  of dimension at most  $k-1$  in which  $\mathbf{T}'(\mathbb{R}^R)$  is included in. Thus, there exists  $\boldsymbol{\theta}^t$  such that  $\boldsymbol{\lambda}(\boldsymbol{\theta}^t) \in \mathbb{R}^k$  is a non-zero vector that is orthogonal to  $\mathbf{T}'(\mathbb{R}^R)$ . Because of the orthogonality, it holds for all  $x_t, \dots, x^{t-R} \in \mathbb{R}$  that  $\sum_{j=1}^k T'_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) = 0$ . By integrating, we find that  $\sum_{j=1}^k T_j(x^t, \{x^{t-1:t-R}\}) \lambda_j(\boldsymbol{\theta}^t) = \text{const}$ . Since  $\boldsymbol{\lambda}(\boldsymbol{\theta}^t) \neq \mathbf{0}$  and  $l(\mathbb{R}) \neq 0$ , the distribution cannot be strongly exponential, which contradicts the hypothesis.

**Lemma 4.** *Consider a strongly exponential autoregressive distribution of size  $k \geq 2$  for which the sufficient statistics  $\mathbf{T}$  are twice differentiable almost everywhere. Then it holds that*

$$\text{rank}((T'_1(x^t, \{x^{t-1:t-R}\}), T''_1(x^t, \{x^{t-1:t-R}\})^\top, \dots, (T'_k(x^t, \{x^{t-1:t-R}\}), T''_k(x^t, \{x^{t-1:t-R}\})^\top) \geq 2 \quad (22)$$

almost everywhere on  $\mathbb{R}$ .

*Proof:* Suppose there exists a set  $\mathcal{X}$  so that  $l(\mathcal{X}) > 0$ , but the equation (22) does not hold. In other words, for all  $j \in \{1, \dots, k\}$  and  $x \in \mathcal{X}$ , the vectors  $(T'_j(x^t, \{x^{t-1:t-R}\}), T''_j(x^t, \{x^{t-1:t-R}\})^\top)$  are collinear. This means that there exists a vector  $\boldsymbol{\alpha} \in \mathbb{R}^k$ ,  $\boldsymbol{\alpha} \neq \mathbf{0}$ , so that  $\sum_{j=1}^k \alpha_j T'_j(x^t, \{x^{t-1:t-R}\}) = 0$ . By integrating, we get  $\sum_{j=1}^k \alpha_j T_j(x^t, \{x^{t-1:t-R}\}) = \text{const}$  for all  $x \in \mathcal{X}$ . Since  $l(\mathcal{X}) > 0$ , this contradicts the hypothesis.

**Lemma 5.** *Consider  $P$  strongly exponential autoregressive distributions of size  $k \geq 2$  for which the sufficient statistics  $\mathbf{T}_j$ ,  $j = 1, \dots, P$  are twice differentiable almost everywhere. Let  $\mathbf{x} := (x_1, \dots, x_P) \in \mathbb{R}^P$  and  $\mathbf{e}^{(j,i)}(x_i) = (0, \dots, 0, T'_{j,i}(x_i), T''_{j,i}(x_i), 0, \dots, 0) \in \mathbb{R}^{2P}$ , so that the non-zero entries are at indices  $(2j, 2j+1)$ . Then the matrix  $\mathbf{E}(\mathbf{z}) = (\mathbf{e}^{(1,1)}(x_1), \dots, \mathbf{e}^{(1,k)}(x_1), \dots, \mathbf{e}^{(P,1)}(x_P), \dots, \mathbf{e}^{(P,k)}(x_P)) \in \mathbb{R}^{2P \times Pk}$  has rank  $2P$  almost everywhere on  $\mathbb{R}^P$ .*

*Proof:* As the non-zero entries are at indices  $(2j, 2j+1)$ , and there are  $k$  columns in the matrix  $\mathbf{E}$  for each  $j = 1, \dots, P$ , the matrix  $\mathbf{E}$  has at least the rank of  $P$ . By using Lemma 4, it can be deduced that for each  $j = 1, \dots, P$ , the submatrix  $\mathbf{E}_j = (\mathbf{e}^{(j,1)}(x_j), \dots, \mathbf{e}^{(j,k)}(x_j))$  has rank greater or equal to 2 almost everywhere on  $\mathbb{R}$ . Thus, it can be concluded that the rank of  $\mathbf{E}$  is  $2P$  almost everywhere on  $\mathbb{R}^P$ . □

## B Proofs

In this section, the proofs are provided for the main identifiability theorems and for all propositions. The proofs of Theorems 1 and 2 closely follow the approach of [1], where the identifiability was proved for the exponential family without autoregressive structure.

### B.1 Proof of Proposition 1

Since  $(\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$  is identifiable up to block-affine transformation, we have  $\tilde{\mathbf{T}}(\tilde{\mathbf{z}}) = \mathbf{A}\mathbf{T}(\mathbf{z}) + \mathbf{c}$ , where  $\mathbf{A}$  is a block-permutation matrix and  $\mathbf{c}$  is a constant vector.

Let  $\pi$  be the permutation of  $\{1, \dots, P\}$  induced by the block structure of  $\mathbf{A}$ . For each  $i$ , the  $i$ th block equation of the above is:  $\tilde{\mathbf{T}}_i(\tilde{z}_i) = \mathbf{A}_{i,\pi(i)}\mathbf{T}_{\pi(i)}(z_{\pi(i)}) + \mathbf{c}_i$  where  $\mathbf{A}_{i,\pi(i)}$  is the  $k \times k$  submatrix of  $\mathbf{A}$  corresponding to the transformation from the  $\pi(i)$ th to the  $i$ th component, and  $\mathbf{c}_i \in \mathbb{R}^k$  is the corresponding subvector of  $\mathbf{c}$ .

By applying  $\tilde{g}_i$  to both sides for each  $i$  and using assumption (ii), we have  $a_i \tilde{z}_i = \tilde{g}_i(\mathbf{A}_{i,\pi(i)}\mathbf{T}_{\pi(i)}(z_{\pi(i)}) + \mathbf{c}_i)$ . Let  $g_{\pi(i)}(z_{\pi(i)}) = \frac{1}{a_i} \tilde{g}_i(\mathbf{A}_{i,\pi(i)}\mathbf{T}_{\pi(i)}(z_{\pi(i)}) + \mathbf{c}_i)$ . Then, we have  $\tilde{z}_i = g_{\pi(i)}(z_{\pi(i)})$ .

The permutation  $\pi$  defines a permutation matrix  $\mathbf{P}$ , giving us  $\tilde{\mathbf{z}} = \mathbf{P}(g_1(z_1), \dots, g_P(z_P))^\top$ .  $\square$

### B.2 Proof of Theorem 1

**Step 1.** Let us denote  $\mathbf{x}^- = \{\mathbf{x}^{t-1:t-R}\}$ ,  $\mathbf{x} = \mathbf{x}^t$  and  $\mathbf{z} = \mathbf{z}^t$ . Suppose there are two sets of parameters  $\boldsymbol{\theta} = (\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$  and  $\tilde{\boldsymbol{\theta}} = (\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}})$  such that  $p_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}}(\mathbf{x}|\mathbf{x}^-, \mathbf{u}) = p_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\mathbf{x}|\mathbf{x}^-, \mathbf{u})$  for all  $(\mathbf{x}|\mathbf{x}^-, \mathbf{u})$ . Then

$$\int_{\mathcal{Z}} p_{\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda}}(\mathbf{x}, \mathbf{z}|\mathbf{x}^-, \mathbf{u}) d\mathbf{z} = \int_{\mathcal{Z}} p_{\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\mathbf{x}, \mathbf{z}|\mathbf{x}^-, \mathbf{u}) d\mathbf{z} \quad (23)$$

$$\implies \int_{\mathcal{Z}} p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{z}|\mathbf{x}^-, \mathbf{u}) p_{\mathbf{f}}(\mathbf{x}|\mathbf{z}) d\mathbf{z} = \int_{\mathcal{Z}} p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\mathbf{z}|\mathbf{x}^-, \mathbf{u}) p_{\tilde{\mathbf{f}}}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad (24)$$

$$\stackrel{(i)}{\implies} \int_{\mathcal{Z}} p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{z}|\mathbf{x}^-, \mathbf{u}) p_{\epsilon}(\mathbf{x} - \mathbf{f}(\mathbf{z})) d\mathbf{z} = \int_{\mathcal{Z}} p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\mathbf{z}|\mathbf{x}^-, \mathbf{u}) p_{\epsilon}(\mathbf{x} - \tilde{\mathbf{f}}(\mathbf{z})) d\mathbf{z} \quad (25)$$

$$\implies \int_{\mathcal{X}} p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{q}(\bar{\mathbf{x}})|\mathbf{x}^-, \mathbf{u}) p_{\epsilon}(\mathbf{x} - \bar{\mathbf{x}}) |\det(J_{\mathbf{q}}(\bar{\mathbf{x}}))| d\bar{\mathbf{x}} = \int_{\mathcal{X}} p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\tilde{\mathbf{q}}(\bar{\mathbf{x}})|\mathbf{x}^-, \mathbf{u}) p_{\epsilon}(\mathbf{x} - \bar{\mathbf{x}}') |\det(J_{\tilde{\mathbf{q}}}(\bar{\mathbf{x}}'))| d\bar{\mathbf{x}} \quad (26)$$

$$\implies \int_{\mathbb{R}^S} \tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-}(\mathbf{q}(\bar{\mathbf{x}})) p_{\epsilon}(\mathbf{x} - \bar{\mathbf{x}}) d\bar{\mathbf{x}} = \int_{\mathbb{R}^S} \tilde{p}_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{f}}, \mathbf{u}, \mathbf{x}^-}(\tilde{\mathbf{q}}(\bar{\mathbf{x}})) p_{\epsilon}(\mathbf{x} - \bar{\mathbf{x}}') d\bar{\mathbf{x}} \quad (27)$$

$$\implies (\tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-} * p_{\epsilon})(\bar{\mathbf{x}}) = (\tilde{p}_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{f}}, \mathbf{u}, \mathbf{x}^-} * p_{\epsilon})(\bar{\mathbf{x}}') \quad (28)$$

$$\implies F[\tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-}](\omega) \varphi_{\epsilon}(\omega) = F[\tilde{p}_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{f}}, \mathbf{u}, \mathbf{x}^-}](\omega) \varphi_{\epsilon}(\omega) \quad (29)$$

$$\stackrel{(i)}{\implies} F[\tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-}](\omega) = F[\tilde{p}_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{f}}, \mathbf{u}, \mathbf{x}^-}](\omega) \quad (30)$$

$$\implies \tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-}(\mathbf{x}) = \tilde{p}_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{f}}, \mathbf{u}, \mathbf{x}^-}(\mathbf{x}) \quad (31)$$

- In equation (26),  $J$  denotes Jacobian, a variable change  $\bar{\mathbf{x}} = \mathbf{f}(\mathbf{z})$  is introduced left hand side and  $\bar{\mathbf{x}}' = \tilde{\mathbf{f}}(\mathbf{z})$  to right hand side.
- In equation (27),  $\tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-} = p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{q}(\bar{\mathbf{x}})|\mathbf{x}^-, \mathbf{u}) |\det(J_{\mathbf{q}}(\bar{\mathbf{x}}))| \mathbb{1}_{\mathcal{X}}(\mathbf{x})$  is introduced left hand side and similarly to right hand side. The indicator function  $\mathbb{1}_{\mathcal{X}}(\mathbf{x})$  is defined as  $\mathbb{1}_{\mathcal{X}}(\mathbf{x}) = \begin{cases} 1, & \text{when } \mathbf{x} \in \mathcal{X}, \\ 0, & \text{otherwise.} \end{cases}$
- In equation (28),  $*$  denotes a convolution operator.
- In equation (29),  $F$  denotes Fourier transform, and  $\varphi_{\epsilon} = F[p_{\epsilon}]$ .
- In equation (30),  $\varphi_{\epsilon}$  is dropped from both sides because of assumption (i) ( $\varphi_{\epsilon}$  is non-zero almost everywhere).

The step 1 guarantees that if the distributions with noise  $\epsilon$  are the same, then the noise-free distributions have to be the same.

**Step 2.** By starting from equation (31) and replacing the conditioning variable  $\mathbf{x}^-$  with  $\mathbf{q}(\mathbf{x}^-) = \{\mathbf{q}_x^{t-1}, \dots, \mathbf{q}_x^{t-R}\}$  (this can be done because  $\mathbf{f}(\mathbf{q}(\mathbf{x})) = \mathbf{x}$ , meaning that  $\mathbf{q}(\mathbf{x})$  contains the same information as  $\mathbf{x}$ ), denoting that the transformation  $\mathbf{q}$  is applied to all  $\mathbf{x}^i$ ,  $i = t-1, \dots, t-R$ , we get the following form:

$$\tilde{p}_{\mathbf{T}, \boldsymbol{\lambda}, \mathbf{f}, \mathbf{u}, \mathbf{x}^-}(\mathbf{x}) = \tilde{p}_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}, \tilde{\mathbf{f}}, \mathbf{u}, \mathbf{x}^-}(\mathbf{x}) \quad (32)$$

$$\implies p_{\mathbf{T}, \boldsymbol{\lambda}}(\mathbf{q}(\mathbf{x}) | \mathbf{q}(\mathbf{x}^-), \mathbf{u}) |\det(J_{\mathbf{q}}(\mathbf{x}))| \mathbb{1}_{\mathcal{X}}(\mathbf{x}) = p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}(\tilde{\mathbf{q}}(\mathbf{x}) | \mathbf{q}(\mathbf{x}^-), \mathbf{u}) |\det(J_{\tilde{\mathbf{q}}}(\mathbf{x}))| \mathbb{1}_{\mathcal{X}}(\mathbf{x}) \quad (33)$$

By taking a logarithm on both sides of equation (33) and replacing  $p_{\mathbf{T}, \boldsymbol{\lambda}}$  and  $p_{\tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}}}$  with the form in equation (5), we get:

$$\begin{aligned} \log |\det(J_{\mathbf{q}}(\mathbf{x}))| + \sum_{i=1}^P (\log Q_i(q_i(\mathbf{x}), q_i(\mathbf{x}^-)) - \log Z_i(\mathbf{u}) + \sum_{j=1}^k T_{i,j}(q_i(\mathbf{x}), q_i(\mathbf{x}^-)) \lambda_{i,j}(\mathbf{u})) = \\ \log |\det(J_{\tilde{\mathbf{q}}}(\mathbf{x}))| + \sum_{i=1}^P (\log \tilde{Q}_i(\tilde{q}_i(\mathbf{x}), \tilde{q}_i(\mathbf{x}^-)) - \log \tilde{Z}_i(\mathbf{u}) + \sum_{j=1}^k \tilde{T}_{i,j}(\tilde{q}_i(\mathbf{x}), \tilde{q}_i(\mathbf{x}^-)) \tilde{\lambda}_{i,j}(\mathbf{u})) \end{aligned} \quad (34)$$

Let  $\mathbf{u}_0, \dots, \mathbf{u}_{Pk}$  be the distinct points in assumption (iv). Then, we have  $Pk + 1$  equations as in (34), one for each point. By subtracting the first equation from the others, for point  $\mathbf{u}_l$ ,  $l = 1, \dots, Pk$ , we have

$$\begin{aligned} \sum_{i=1}^P \log \frac{Z_i(\mathbf{u}_0)}{Z_i(\mathbf{u}_l)} + \sum_{i=1}^P \sum_{j=1}^k (T_{i,j}(q_i(\mathbf{x}), q_i(\mathbf{x}^-)) (\lambda_{i,j}(\mathbf{u}_l) - \lambda_{i,j}(\mathbf{u}_0))) = \\ \sum_{i=1}^P \log \frac{\tilde{Z}_i(\mathbf{u}_0)}{\tilde{Z}_i(\mathbf{u}_l)} + \sum_{i=1}^P \sum_{j=1}^k (\tilde{T}_{i,j}(\tilde{q}_i(\mathbf{x}), \tilde{q}_i(\mathbf{x}^-)) (\tilde{\lambda}_{i,j}(\mathbf{u}_l) - \tilde{\lambda}_{i,j}(\mathbf{u}_0))) \end{aligned} \quad (35)$$

Let us define  $\bar{\boldsymbol{\lambda}}(\mathbf{u}) = \boldsymbol{\lambda}(\mathbf{u}) - \boldsymbol{\lambda}(\mathbf{u}_0)$ , and subtract  $\sum_{i=1}^P \log \frac{\tilde{Z}_i(\mathbf{u}_0)}{\tilde{Z}_i(\mathbf{u}_l)}$  both sides. Then we have

$$\sum_{i=1}^P \sum_{j=1}^k (T_{i,j}(q_i(\mathbf{x}), q_i(\mathbf{x}^-)) (\bar{\lambda}_{i,j}(\mathbf{u}_l))) = \sum_{i=1}^P \log \frac{Z_i(\mathbf{u}_0) \tilde{Z}_i(\mathbf{u}_0)}{Z_i(\mathbf{u}_l) \tilde{Z}_i(\mathbf{u}_l)} + \sum_{i=1}^P \sum_{j=1}^k (\tilde{T}_{i,j}(\tilde{q}_i(\mathbf{x}), \tilde{q}_i(\mathbf{x}^-)) (\bar{\lambda}_{i,j}(\mathbf{u}_l))) \quad (36)$$

Let us write  $b_l = \sum_{i=1}^P \log \frac{Z_i(\mathbf{u}_0) \tilde{Z}_i(\mathbf{u}_0)}{Z_i(\mathbf{u}_l) \tilde{Z}_i(\mathbf{u}_l)}$  and set  $\mathbf{b} = (b_1, \dots, b_{Pk})$ . Let  $\mathbf{L}$  be the matrix in assumption (iv), and  $\tilde{\mathbf{L}}$  similar matrix for  $\tilde{\boldsymbol{\lambda}}$ . By expressing (36) in matrix form for all point  $b_l$ ,  $l = 1, \dots, Pk$ , we have:

$$\mathbf{L}^\top \mathbf{T}(\mathbf{q}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) = \tilde{\mathbf{L}}^\top \tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) + \mathbf{b} \quad (37)$$

$$\implies \mathbf{T}(\mathbf{q}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) = (\mathbf{L}^\top)^{-1} \tilde{\mathbf{L}}^\top \tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) + (\mathbf{L}^\top)^{-1} \mathbf{b} \quad (38)$$

$$\implies \mathbf{T}(\mathbf{q}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) = \mathbf{A} \tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) + \mathbf{c}, \quad (39)$$

where  $\mathbf{A} = (\mathbf{L}^\top)^{-1} \tilde{\mathbf{L}}^\top$  and  $\mathbf{c} = (\mathbf{L}^\top)^{-1} \mathbf{b}$ .

**Step 3.** By assumption (iii), Jacobian of  $\mathbf{T}$  exists and is a  $Pk \times P$  matrix of rank  $P$ . Because equation (39) holds, it also holds that  $J(\mathbf{T}(\mathbf{q}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-))) = \mathbf{A} J(\tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \tilde{\mathbf{q}}(\mathbf{x}^-)))$  and that  $\text{rank}(J(\mathbf{T}(\mathbf{q}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)))) = \text{rank}(J(\tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \tilde{\mathbf{q}}(\mathbf{x}^-))))$ . This leads to the fact that both  $\mathbf{A}$  and  $J(\tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \tilde{\mathbf{q}}(\mathbf{x}^-)))$  are of rank  $P$ .

- If  $k = 1$ , then  $\mathbf{A}$  is invertible since it is a  $P \times P$  matrix of rank  $P$ .
- If  $k \geq 2$ , define  $\bar{\mathbf{z}} = \mathbf{q}(\mathbf{x})$ ,  $\bar{\mathbf{z}}^- = \mathbf{q}(\mathbf{x}^-)$  and  $\mathbf{T}_i = (T_{i,1}(\bar{z}_i, \bar{z}_i^-), \dots, T_{i,k}(\bar{z}_i, \bar{z}_i^-))$ . Based on Lemma 3, it holds that for each  $i = 1, \dots, P$ , there exists  $k$  points  $(\bar{z}_i^j, \bar{z}_i^{-,j})$ ,  $j = 1, \dots, k$  such that  $(\mathbf{T}_i'(\bar{z}_i^1, \bar{z}_i^{-,1}), \dots, \mathbf{T}_i'(\bar{z}_i^k, \bar{z}_i^{-,k}))$  are linearly independent. Let us define  $\mathbf{Q} = (J(\mathbf{T}(\bar{\mathbf{z}}^1, \bar{\mathbf{z}}^{-,1})), \dots, J(\mathbf{T}(\bar{\mathbf{z}}^k, \bar{\mathbf{z}}^{-,k})))$ , where each Jacobian is  $Pk \times P$  matrix calculated with respect to  $\bar{\mathbf{z}}^i$ , and the vector  $\bar{\mathbf{z}}^l$  and  $\bar{\mathbf{z}}^{-,l}$  are defined as  $\bar{\mathbf{z}}^l = (\bar{z}_1^l, \dots, \bar{z}_P^l)$  and  $\bar{\mathbf{z}}^{-,l} = (\bar{z}_1^{-,l}, \dots, \bar{z}_P^{-,l})$ . Similarly define matrix  $\tilde{\mathbf{Q}}$  for Jacobians of  $\tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{f}(\bar{\mathbf{x}}^l)), \tilde{\mathbf{q}}(\mathbf{f}(\bar{\mathbf{x}}^{-,l})))$  for the same points  $l = 1, \dots, k$ . Then, by differentiating the equation (39) for each  $\mathbf{x}_l$ , we get the following in matrix form:

$$\mathbf{Q} = \mathbf{A} \tilde{\mathbf{Q}}. \quad (40)$$

The matrix  $\mathbf{Q}$  is invertible based on Lemma 3, and hence also  $\mathbf{A}$  and  $\tilde{\mathbf{Q}}$  are invertible. As we have invertible  $\mathbf{A}$ , the equation (39) says that the sufficient statistics are identifiable up to linear transformation and a constant.  $\square$

### B.3 Proof of Theorem 2

**Step 1.** The assumptions of theorem 1 holds, so we have

$$\mathbf{T}(\mathbf{q}(\mathbf{x}), \mathbf{q}(\mathbf{x}^-)) = \mathbf{A}\tilde{\mathbf{T}}(\tilde{\mathbf{q}}(\mathbf{x}), \tilde{\mathbf{q}}(\mathbf{x}^-)) + \mathbf{c}, \quad (41)$$

where  $\mathbf{c}$  is a constant vector and  $\mathbf{A}$  is an invertible  $Pk \times Pk$  matrix. Let  $(i, l, a, b)$  be four indices so that  $1 \leq i \leq P$ ,  $1 \leq l \leq k$  refer to the rows of the matrix  $\mathbf{A}$ , and  $1 \leq a \leq P$ ,  $1 \leq b \leq k$  refer to the columns of  $\mathbf{A}$ . Let  $\mathbf{v}(\mathbf{z}) = \tilde{\mathbf{q}}(\mathbf{f}(\mathbf{z})) : \mathcal{Z} \rightarrow \mathcal{Z}$ . The function  $\mathbf{v}$  is bijective as  $\tilde{\mathbf{f}} : \mathcal{Z} \rightarrow \mathcal{X}$  and  $\mathbf{f} : \mathcal{Z} \rightarrow \mathcal{X}$  are injective functions, and  $\tilde{\mathbf{q}}(\mathbf{f}(\mathbf{z})) = \mathbf{z}$ . Further, let there be two other indices  $c, d \in \{1, \dots, P\}$ ,  $c < d$  and denote  $v_i^c = \frac{\partial v_i}{\partial v_c}$  and  $v_i^{c,d} = \frac{\partial v_i}{\partial v_c \partial v_d}$ . By differentiating (41) with respect to  $z_c$ , we get for each  $1 \leq i \leq P$  and  $1 \leq l \leq k$  the following:

$$\frac{\partial T_{i,l}(z_i, z_i^-)}{\partial z_c} = \sum_{a,b} A_{i,l,a,b} \left( \frac{\partial \tilde{T}_{a,b}(v_a(\mathbf{z}), v_a(\mathbf{z}^-))}{\partial v_a(\mathbf{z})} \frac{\partial v_a(\mathbf{z})}{\partial z_c} + \sum_{r=1}^R \frac{\partial \tilde{T}_{a,b}(v_a(\mathbf{z}), v_a(\mathbf{z}^-))}{\partial v_a(\mathbf{z}^{-r})} \frac{\partial v_a(\mathbf{z}^{-r})}{\partial z_c} \right). \quad (42)$$

It holds that  $\frac{\partial v_a(\mathbf{z}^{-r})}{\partial z_c} = 0$  for all  $r = 1, \dots, R$ , as the values of previous time points do not depend on the value of current time point. Thus, we have

$$\frac{\partial T_{i,l}(z_i, z_i^-)}{\partial z_c} = \sum_{a,b} A_{i,l,a,b} \left( \frac{\tilde{T}_{a,b}}{\partial v_a(\mathbf{z})} \frac{\partial v_a(\mathbf{z})}{\partial z_c} \right). \quad (43)$$

By differentiating (43) with respect to  $z_d$ , we get

$$0 = \sum_{a,b} A_{i,l,a,b} \left( \frac{\partial \tilde{T}_{a,b}(v_a(\mathbf{z}), v_a(\mathbf{z}^-))}{\partial z_d} \frac{\partial v_a(\mathbf{z})}{\partial z_c \partial z_d} + \frac{\partial \tilde{T}_{a,b}(v_a(\mathbf{z}))}{\partial^2 v_a(\mathbf{z})} \frac{\partial v_a(\mathbf{z})}{\partial z_c} \frac{\partial v_a(\mathbf{z})}{\partial z_d} \right). \quad (44)$$

Let us define  $\mathbf{r}_a^1(\mathbf{z}) = (v_a^{1,2}(\mathbf{z}), \dots, v_a^{P-1,P}(\mathbf{z})) \in \mathbb{R}^{\frac{P(P-1)}{2}}$ ,  $\mathbf{r}_a^2(\mathbf{z}) = (v_a^1(\mathbf{z})v_a^2(\mathbf{z}), \dots, v_a^{P-1}(\mathbf{z})v_a^P(\mathbf{z})) \in \mathbb{R}^{\frac{P(P-1)}{2}}$ ,  $\mathbf{M}(\mathbf{z}) = (\mathbf{r}_1^1(\mathbf{z}), \mathbf{r}_1^2(\mathbf{z}), \dots, \mathbf{r}_P^1(\mathbf{z}), \mathbf{r}_P^2(\mathbf{z})) \in \mathbb{R}^{\frac{P(P-1)}{2} \times \frac{P(P-1)}{2}}$  and  $\mathbf{e}^{(a,b)}(z_i) = (0, \dots, 0, T'_{a,b}(z_i), T''_{a,b}(z_i), 0, \dots, 0) \in \mathbb{R}^{2P}$ , so that the non-zero entries are at indices  $(2a, 2a+1)$  and  $\mathbf{E}(\mathbf{z}) = (\mathbf{e}^{(1,1)}(z_1), \dots, \mathbf{e}^{(1,k)}(z_1), \dots, \mathbf{e}^{(P,1)}(z_P), \dots, \mathbf{e}^{(P,k)}(z_P)) \in \mathbb{R}^{2P \times Pk}$ . Finally, let  $A_{i,l}$  be the  $(i, l)$ th row of the matrix  $\mathbf{A}$ . Then, by gathering the equation (43) for all pairs  $(c, d)$ ,  $c < d$  and pairs  $(i, l)$  to a matrix form, we get

$$\mathbf{M}(\mathbf{z})\mathbf{E}(\mathbf{z})\mathbf{A} = \mathbf{0}. \quad (45)$$

By Lemma 5, the matrix  $\mathbf{E}$  is of rank  $2P$  almost surely on  $\mathcal{Z}$ . Since the matrix  $\mathbf{A}$  is full rank  $Pk \times Pk$  matrix, we have  $\text{rank}(\mathbf{EA}) = 2P$  almost surely on  $\mathcal{Z}$ . Hence, by multiplying (45) from right with the pseudo-inverse of  $(\mathbf{EA})$  we have

$$\mathbf{M}(\mathbf{z}) = \mathbf{0}. \quad (46)$$

Particularly,  $\mathbf{r}_a^2 = \mathbf{0}$  for all  $a = 1, \dots, P$ . This means that at each  $\mathbf{z} \in \mathcal{Z}$ , the Jacobian of  $\mathbf{v}$ ,  $J_{\mathbf{v}}$  has at most one non-zero entry in each row. Because  $J_{\mathbf{v}}$  is invertible and continuous, the locations of the non-zero entries are fixed and do not change as function of  $\mathbf{z}$ . This proves that the function  $\tilde{\mathbf{q}}(\mathbf{f}(\mathbf{z}))$  is a composition of a permutation and a point-wise nonlinearity.

**Step 2.** Without loss of generality, we assume that the permutation in  $\mathbf{v}$  is identity. Let  $\bar{\mathbf{T}}(\mathbf{z}) = \tilde{\mathbf{T}}(\mathbf{v}(\mathbf{z})) + \mathbf{A}^{-1}\mathbf{c}$ . In particular,  $\bar{\mathbf{T}}$  is then a point-wise nonlinearity. Then, the equation (41) can be written as

$$\mathbf{T}(\mathbf{z}, \mathbf{z}^-) = \mathbf{A}\bar{\mathbf{T}}(\mathbf{z}, \mathbf{z}^-). \quad (47)$$

Let  $\mathbf{W} = \mathbf{A}^{-1}$ . Then, the equation (47) can be written for each component  $1 \leq i \leq P$  and sufficient statistic  $1 \leq l \leq k$  as

$$\bar{T}_{i,l} = \sum_{a,b} D_{i,l,a,b} T_{a,b}(z_a, z_a^-). \quad (48)$$

By differentiating both sides with respect to  $z_c$ ,  $c \neq i$ , we get

$$0 = \sum_b D_{i,l,c,b} \frac{\partial T'_{c,b}(z_a, z_a^-)}{\partial z_c}. \quad (49)$$

By Lemma 1, we know that  $D_{i,l,c,b} = 0$  for all  $1 \leq b \leq k$ , and since (49) holds for all  $l$  and  $c \neq i$ , the matrix  $\mathbf{D}$  must have a block diagonal form

$$\mathbf{D} = \begin{pmatrix} \mathbf{D}_1 & & \\ & \ddots & \\ & & \mathbf{D}_P \end{pmatrix}, \quad (50)$$

where each submatrix  $\mathbf{D}_1, \dots, \mathbf{D}_P$  is a  $k \times k$  matrix. Then, also the matrix  $\mathbf{A}$  has the same block diagonal form, meaning that each submatrix  $\mathbf{A}_i$  transforms  $\mathbf{T}_i(\mathbf{z}, \mathbf{z}^-)$  into  $\bar{\mathbf{T}}_i(\mathbf{z}, \mathbf{z}^-)$ . Since  $\bar{\mathbf{T}}$  is a point-wise nonlinearity,  $\mathbf{A}$  has to be a permutation matrix.  $\square$

#### B.4 Proof of Proposition 2

Based on the assumptions we have the following equalities

$$\begin{aligned} \tilde{z}_j &= a_{11}z_i + a_{12}z_i^2 + c_1, \\ \tilde{z}_j^2 &= a_{21}z_i + a_{22}z_i^2 + c_2, \end{aligned}$$

for some constants  $a_{11}, a_{12}, a_{21}, a_{22}, c_1$  and  $c_2$ . By squaring the first equation, we have  $(a_{11}z_i + a_{12}z_i^2 + c_1)^2 = a_{21}z_i + a_{22}z_i^2 + c_2$ . In order the equation to hold for all  $z_i \in \mathcal{Z}$ , it must hold that  $a_{12} = 0$ . Hence, we have that  $\tilde{z}_j = a_{11}z_i + c_1$ .  $\square$

#### B.5 Proof of Proposition 3

Since  $(\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$  is identifiable up to block-affine transformation, we have  $\tilde{\mathbf{T}}(\tilde{\mathbf{z}}) = \mathbf{A}\mathbf{T}(\mathbf{z}) + \mathbf{c}$ , where  $\mathbf{A}$  is a block-permutation matrix and  $\mathbf{c}$  is a constant vector.

Let  $\pi$  be the permutation of  $\{1, \dots, P\}$  induced by the block structure of  $\mathbf{A}$ , and  $j = \pi_i$ . Then we have that  $\tilde{\mathbf{T}}_j(\tilde{z}_j^t, \dots, \tilde{z}_j^{t-R}) = \mathbf{A}_{i,j} \mathbf{T}_i(z_j^t, \dots, z_j^{t-R})$ , where  $\mathbf{A}_{i,j}$  is a  $k \times k$  submatrix of  $\mathbf{A}$  corresponding the indices  $i$  and  $j$ . Because of Gaussian AR form (1), we have

$$\begin{aligned} p(\mathbf{z} | \{\mathbf{z}^{t-1:t-R}\}, \mathbf{u}^t, \dots, \mathbf{u}^{t-R}) = \\ \prod_{i=1}^P \frac{1}{2\pi\sigma_i(\mathbf{u}^t)} \exp \left[ -\frac{\left( z_i - \mu_i(\mathbf{u}^t) - \sum_{r=1}^R (\gamma_r(\mathbf{u}^t) z_i^{t-r} - \mu_i(\mathbf{u}^{t-r})) \right)^2}{2\sigma^2(\mathbf{u}^t)} \right]. \end{aligned} \quad (51)$$

and similar form for  $\tilde{z}_j$  with parameter functions  $\tilde{\mu}_j, \tilde{\sigma}_j, \tilde{\gamma}_{j,1}, \dots, \tilde{\gamma}_{j,R}$ . Let  $\gamma_{i,r} := \gamma_{i,r}(\mathbf{u}^t)$ ,  $\mu_{i,r} := \mu_i(\mathbf{u}^{t-r})$  and  $\sigma_i := \sigma_i(\mathbf{u}^t)$ . By expanding the nominator in the exponential term, we have

$$\begin{aligned} (z_i^t)^2 - 2z_i^t \mu_{i,0} - 2z_i^t \sum_{r=1}^R \gamma_{i,r} z_i^{t-r} + 2z_i^t \sum_{r=1}^R \gamma_{i,r} \mu_{i,r} + \mu_{i,0}^2 + 2\mu_{i,0} \sum_{r=1}^R \gamma_{i,r} z_i^{t-r} + \\ \left( \sum_{r=1}^R \gamma_{i,r} z_i^{t-r} \right)^2 - 2 \left( \sum_{r=1}^R \gamma_{i,r} z_i^{t-r} \right) \left( \sum_{r=1}^R \gamma_{i,r} \mu_{i,r} \right) + \left( \sum_{r=1}^R \gamma_{i,r} \mu_{i,r} \right)^2. \end{aligned} \quad (52)$$

From this form, it is easy to see that the minimal sufficient statistics are  $T_{i,1} = (z_i^t)^2$ ,  $T_{i,2} = z_i^t$ ,  $T_{i,3,r} = z_i^t z_i^{t-r}$ ,  $T_{i,4,r} = z_i^{t-r}$ ,  $T_{i,5,r_1,r_2} = z_i^{t-r_1} z_i^{t-r_2}$ ,  $r, r_1, r_2 \in \{1, \dots, R\}$ . Similarly, we have the sufficient statistics  $\mathbf{T}_j(z_j^t, \dots, z_j^{t-R})$ . Because of the block-affine identifiability, we have for each  $k_1 \in \{1, \dots, k\}$  that

$$\tilde{T}_{k_1,j} = \sum_{k_2=1}^k a_{k_2,k_1,i} T_{k_2,i} + c_{i,k_1}, \quad (53)$$

where  $a_{k_1,k_2,i}$  and  $c_{i,k_1}$  are constants. Importantly, we have for all  $r = 0, \dots, R$  that  $\tilde{z}_j^{t-r} = \sum_{k_2=1}^k a_{k_2,r_1,i} T_{k_2,i} + c_i$  and  $(\tilde{z}_j^{t-r})^2 = \sum_{k_2=1}^k a_{k_2,r_2,i} T_{k_2,i} + c_i$ . By squaring the first equation, we have that

$$\left( \sum_{k_2=1}^k a_{k_2,r_1,i} T_{k_2,i} + c_{i,r_1} \right)^2 = \sum_{k_2=1}^k a_{k_2,r_2,i} T_{k_2,i} + c_{i,r_2}. \quad (54)$$

This equation holds only if the coefficients of the third order and above in the left hand side are zero, meaning that  $a_{1,r_1,i}, a_{(3,r),r_1,i}, a_{(5,r),r_1,i} = 0$ . Hence, we have for all  $r_1 = 0, \dots, R$  and  $t = R+1, \dots, T$  that

$$\tilde{z}_j^{t-r_1} = \sum_{r_2=0}^R b_{r_1,r_2,i} z_i^{t-r_2} + c_{r_1,i}, \quad (55)$$

where  $b_{r_1,r_2,i}$  are constants. Since (55) holds for all  $t = R+1, \dots, T$ , we also have the following equations:

$$\begin{aligned} \tilde{z}_j^t &= \sum_{r=0}^R b_{0,r,i} z_i^{t-r} + c_{0,i}, \\ \tilde{z}_j^t &= \sum_{r=0}^R b_{R,r,i} z_i^{t+R-r} + c_{R,i}, \end{aligned} \quad (56)$$

where the second equation is obtained by shifting (55), for  $r_1 = R$ ,  $R$  time steps forward. From (56) we can deduce that all coefficients  $b_{0,r,i}$ ,  $r \neq 0$ , have to be zero in order for the equations to hold for all  $t \in \{R+1, T\}$ . Hence, we obtain  $\tilde{z}_j^t = b_{0,0,i} z_i^t + c_{0,i}$ , which concludes the proof.  $\square$

## B.6 Proof of Theorem 3

The lower bound of the data log likelihood (ELBO) (9) can also be written in the following format:

$$\text{ELBO} = E_{q_{\theta}(\mathbf{z}|\mathbf{x},\mathbf{u})}(\log p_{\theta}(\mathbf{x}|\mathbf{x}^-, \mathbf{u}) + \text{KL}(\log q_{\theta_g}(\mathbf{z}|\mathbf{x}, \mathbf{x}^-, \mathbf{u}) || p_{\theta}(\mathbf{z}|\mathbf{x}, \mathbf{x}^-, \mathbf{u}))), \quad (57)$$

where KL is Kullback-Leibler divergence and the set  $(\tilde{\mathbf{f}}, \tilde{\mathbf{T}}, \tilde{\boldsymbol{\lambda}})$  are parametrized by  $\boldsymbol{\theta}$ . Minimizing ELBO given in (9) with respect to the parameters  $(\boldsymbol{\theta}, \boldsymbol{\theta}_g)$  is equivalent to minimizing (57), which means that in the limit of infinite data, the KL term eventually reaches zero, making the loss equal to the data log likelihood. Hence in this case, minimizing ELBO is equivalent to maximum likelihood estimation (MLE). As we assume that Theorem 1 or Theorem 2 hold, the consistency of MLE guarantees that the estimation converges to the corresponding identifiability class of the true set  $(\mathbf{f}, \mathbf{T}, \boldsymbol{\lambda})$  in the limit of infinite data.  $\square$

## C Illustration of iVAEar framework

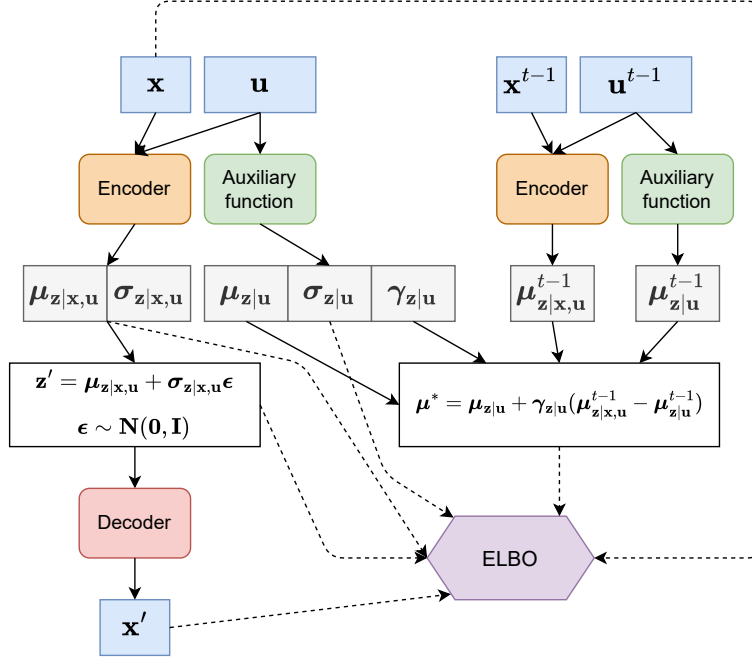
An illustration of iVAEar framework is provided in  $R = 1$  case in Figure 1.

## D Additional simulation details

The parameters used in all simulation settings of Section 4.1, are provided in Table 1.

## References

1. I. Khemakhem, D. Kingma, R. Monti, and A. Hyvärinen, “Variational autoencoders and nonlinear ICA: A unifying framework,” in *International Conference on Artificial Intelligence and Statistics*, pp. 2207–2217, PMLR, 2020.



**Fig. 1.** Schematic presentation of iVAEar algorithm in  $R = 1$  case.

**Table 1.** The parameters for the Matern covariance function in all simulation settings.

	IC1	IC2	IC3	IC4	IC5	IC6
$\phi$	0.20	0.15	0.10	0.30	0.05	0.25
$\nu$	0.50	1.00	0.25	2.00	0.75	1.50