
Multi-scale Graphical Models for Spatio-Temporal Processes: Technical Supplement

Firdaus Janoos* Huseyin Denli Niranjan Subrahmanya
ExxonMobil Corporate Strategic Research
Annandale, NJ 08801

A Proof for Theorem 1

Proof. The original model is given by:

$$\mathbf{x}[t] = \sum_{p=1}^P \mathbf{A}[p] \mathbf{x}[t-p] + \mathbf{u}[t], \quad (1)$$

$$\mathbf{y}[t] = \sum_{q=1}^Q \mathbf{B}[q] \mathbf{y}[t-q] + \mathbf{Z} \mathbf{x}[t] + \mathbf{v}[t]. \quad (2)$$

Define the short-hand notations $\mathbb{E} \{ \cdot \mid t \} \triangleq \mathbb{E} \{ \cdot \mid \mathbf{y}_p; s \leq t \}$ and $\hat{\mathbf{B}}_q \triangleq (\mathbf{Z} \mathbf{Z}^\top) \odot \mathbf{B}_p$, captures the constraint that site $i \leftrightarrow$ site j if they belong to the same global component k . Also to reduce clutter, we use subscripts for time indexing. Suitably adjusting the time-indices, the model of eqn. (1) can be re-written as:

$$\mathbf{x}_{t+1} = \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{x}_{t-p} + \mathbf{u}_t \quad (3)$$

$$\mathbf{y}_t = \sum_{q=1}^Q \hat{\mathbf{B}}_q \mathbf{y}_{t-q} + \mathbf{Z} \mathbf{x}_t + \mathbf{v}_t, \quad (4)$$

It is assumed that the model is stable, that is all the roots of $\det \left| \mathbf{I} - \sum_p \mathbf{A}_p z^{-p} \right|$ (*i.e.* z -transform of \mathbf{A}) lie within the unit circle $|z| < 1$ system. In order to prove Theorem 1 we first introduce the following proposition:

Proposition 1. *Assuming that the system is stable, , then eqn. (3) conditioned on data is equivalent to:*

$$\mathbf{x}_{t+1} \mid t = \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{x}_{t-p} \mid t-p-1 + \sum_{p=0}^{P-1} \mathbf{G}_p \epsilon_{t-p} \quad (5)$$

$$\mathbf{y}_t = \sum_{q=1}^Q \hat{\mathbf{B}}_q \mathbf{y}_{t-q} + \mathbf{Z} \mathbf{x}_t \mid t-1 + \epsilon_t, \quad (6)$$

where $\mathbf{x}_t \mid t-1 = \mathbb{E} \{ \mathbf{x}_t \mid t-1 \}$, $\mathbf{y}_t \mid t-1 = \mathbb{E} \{ \mathbf{y}_t \mid t-1 \}$, $\epsilon_t = y_t - \mathbf{y}_t \mid t-1$, and \mathbf{G} are $K \times N$ matrices defined as follows:

$$\mathbf{G}_p = \sum_{s=p}^{P-1} \mathbf{A}_s \text{Cov} \{ \mathbf{x}_{t-s} \epsilon_t^\top \} \Sigma_\epsilon^{-1} \quad \text{for } p = 0 \dots P-1$$

*Corresponding Author. firdaus@ieee.org

Proof for this proposition is in Appendix A.1. Defining $\hat{\mathbf{X}}, \mathbf{Y}$ and Υ as the z -transform of $\hat{\mathbf{x}}, \mathbf{y}$ and ϵ respectively, we get the z -transform of eqn. (5):

$$\left(\mathbf{I} - \sum_{q=0}^{Q-1} \hat{\mathbf{B}}_q z^{-q} \right) \mathbf{Y} = \left(\mathbf{I} + \mathbf{Z} \left(\mathbf{I} - \sum_{p=0}^{P-1} \mathbf{A}_p z^{-p} \right)^{-1} \left(\sum_{p=0}^{P-1} \mathbf{G}_p \mathbf{z}^{-p} \right) \right) \Upsilon$$

As, \mathbf{Z} is a $N \times K$ matrix of rank K the left pseudo-inverse of \mathbf{Z} is well defined and :

$$\begin{aligned} \left(\mathbf{I} - \sum_{p=0}^{P-1} \mathbf{A}_p z^{-p} \right) \mathbf{Z}^+ \left(\mathbf{I} - \sum_{q=0}^{Q-1} \hat{\mathbf{B}}_q z^{-q} \right) \mathbf{Y} &= \left(\sum_{p=0}^{P-1} \mathbf{G}_p \mathbf{z}^{-p} \right) \Upsilon + \left(\mathbf{I} - \sum_{p=0}^{P-1} \mathbf{A}_p z^{-p} \right) \mathbf{Z}^+ \Upsilon \\ \text{therefore} \\ \mathbf{Z}^* - \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{Z}^+ z^{-p} - \sum_{q=0}^{Q-1} \mathbf{Z}^+ \hat{\mathbf{B}}_q z^{-q} + \sum_{p=0}^{P-1} \sum_{q=0}^{Q-1} \mathbf{A}_p \mathbf{Z}^+ \hat{\mathbf{B}}_q z^{-(p+q)} \\ &= \left(\sum_{p=0}^{P-1} \mathbf{G}_p \mathbf{z}^{-p} \right) \Upsilon + \left(\mathbf{I} - \sum_{p=0}^{P-1} \mathbf{A}_p z^{-p} \right) \mathbf{Z}^+ \Upsilon \end{aligned}$$

Assuming that the system is minimal[3] that is there is strictly no smaller equivalent model, this represents an $N \times N$ rank- K system of auto-regressive order $P + Q$ and moving-average order P . \square

A.1 Proof for Proposition 1

Proof. Firstly, because the assumption on the roots of $\det \left| \mathbf{I} - \sum_p \mathbf{A}_p z^p \right|$, the system is wide sense stationary (WSS) [1]. Now, by the projection property of linear systems the residual $\epsilon_t \perp \mathbf{y}_p; s < t$ and therefore:

$$\mathbb{E} \{ \cdot \mid t \} = \mathbb{E} \{ \cdot \mid t-1 \} + \mathbb{E} \{ \cdot \mid \epsilon_t \}$$

Therefore, the conditional estimate of \mathbf{x}_t :

$$\begin{aligned} \mathbf{x}_{t+1} \mid t &= \sum_{p=0}^{P-1} \mathbb{E} \{ \mathbf{A}_p \mathbf{x}_{t-p} + \mathbf{u}_t \mid t \} \\ &= \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{x}_{t-p} \mid t, \end{aligned}$$

using the property that \mathbf{u}_t independent on $\mathbf{y}_0 \dots \mathbf{y}_t$. Defining $\mathbf{H}_p = \text{Cov} \{ \mathbf{x}_{t-p} \epsilon_t^\top \} \Sigma_\epsilon^{-1}$ and because of the WSS condition, it does not depend on t . Therefore,

$$\begin{aligned} \mathbf{x}_{t+1} \mid t &= \sum_{p=0}^{P-1} \mathbf{A}_p \left(\mathbf{x}_{t-p} \mid t-1 + \mathbb{E} \{ \mathbf{x}_{t-p} \mid \epsilon_t \} \right) \\ &= \mathbf{A}_0 \mathbf{x}_t \mid t-1 + \sum_{p=1}^{P-1} \mathbf{A}_p \mathbf{x}_{t-p} \mid t-1 + \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{H}_p \epsilon_t \\ &= \mathbf{A}_0 \mathbf{x}_t \mid t-1 + \sum_{p=1}^P \mathbf{A}_p \left(\mathbf{x}_{t-p} \mid t-2 + \mathbb{E} \{ \mathbf{x}_{t-p} \mid \epsilon_{t-1} \} \right) + \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{H}_p \epsilon_t \\ &= \mathbf{A}_0 \mathbf{x}_t \mid t-1 + \mathbf{A}_1 \mathbf{x}_{t-1} \mid t-2 + \sum_{p=2}^P \mathbf{A}_p \mathbf{x}_{t-p} \mid t-2 + \sum_{p=1}^P \mathbf{A}_p \mathbf{H}_{p-1} \epsilon_{t-1} + \sum_{p=0}^{P-1} \mathbf{A}_p \mathbf{H}_p \epsilon_t \end{aligned}$$

Continuing this expansion, and setting $G_p = \sum_{s=p}^{P-1} A_s H_{s-p}$ for $p = 0 \dots P-1$, we get:

$$\begin{aligned} \mathbf{x}_{t+1} | t &= \sum_{p=0}^{P-1} A_p \mathbf{x}_{t-p} | t-p-1 + \sum_{p=0}^{P-1} \left[\sum_{s=p}^{P-1} A_s H_{s-p} \right] \epsilon_{t-p} \\ &= \sum_{p=0}^{P-1} A_p \mathbf{x}_{t-p} | t-p-1 + \sum_{p=0}^{P-1} G_p \epsilon_{t-p} \end{aligned} \quad (7)$$

Moreover, since:

$$\mathbf{y}_t | t-1 = \mathbb{E} \left\{ \sum_{q=1}^Q \hat{B}_q \mathbf{y}_{t-q} + Z \mathbf{x}_t + \mathbf{v}_t | t-1 \right\} = \sum_{q=1}^Q \hat{B}_q \mathbf{y}_{t-q} + Z \mathbf{x}_t | t-1$$

we get the following innovations model:

$$\mathbf{y}_t = \sum_{q=1}^Q \hat{B}_q \mathbf{y}_{t-q} + Z \mathbf{x}_t | t-1 + \epsilon_t. \quad (8)$$

Defining $\eta_t = \mathbf{x}_t - \mathbf{x}_t | t-1$, we see that:

$$\epsilon_t = y_t - \mathbf{y}_t | t-1 = Z \mathbf{x}_t + \mathbf{v}_t - Z \mathbf{x}_t | t-1 = Z \eta_t + \mathbf{v}_t \quad (9)$$

while:

$$\begin{aligned} \eta_{t+1} &= \sum_{p=0}^{P-1} A_p (\mathbf{x}_{t-p} - \mathbf{x}_{t-p} | t-1) + \mathbf{u}_t - \sum_{p=0}^{P-1} G_p \epsilon_{t-p} \\ &= \sum_{p=0}^{P-1} A_p \eta_{t-p} + \mathbf{u}_t - \sum_{p=0}^{P-1} G_p (Z \eta_{t-p} + \mathbf{v}_{t-p}) \\ &= \sum_{p=0}^{P-1} (A_p - G_p Z) \eta_{t-p} + \mathbf{u}_t - \sum_{p=0}^{P-1} G_p \mathbf{v}_{t-p} \end{aligned}$$

□

B Proof for Theorem 2

Proof. In order to prove the theorem, we introduce the following proposition:

Proposition 2. *For a 1-d C-D system with infinite boundary conditions and constant Péclet number, for an impulse at the origin $x = 0$, the response in the near-field (i.e. at a point close to the origin) can be approximated by $G(t) \approx \delta(t)$, the Dirac delta function, while in far-field (i.e. at a point far from the origin) and be approximated by a Gaussian function: $G(t) \approx \exp \{-0.5(t - \mu^2 \sigma^{-2})\} / \sqrt{2\pi \sigma^2}$ up to a multiplicative factor of order $\exp \{-\mathcal{O}(t^3)\}$, where μ is equal to the distance and σ^2 is proportional to the product of the distance and the Péclet number.*

The proof to Proposition 2 is given Appendix B.1. Therefore, for the system of Theorem 2, the response at location \mathbf{x}_i to an impulse at $\mathbf{x}_j = 0$ can be approximated by a Gaussian function: $G(t)_{i,j} = \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2) = \exp \{-0.5(t - \mu_{i,j})^2 \sigma_{i,j}^{-2}\} / \sqrt{2\pi \sigma_{i,j}^2}$, where $\mu_{i,j}$ is equal to the distance between \mathbf{x}_i and \mathbf{x}_j and $\sigma_{i,j}$ is proportional to the production of the distance and the Péclet number. Moreover, the response at \mathbf{x}_i to an impulse at \mathbf{x}_i can be approximated by a Dirac delta. Therefore, in Fourier domain the *transfer function* of the 2×2 system consisting of only nodes \mathbf{x}_i and \mathbf{x}_j :

$$\hat{\Psi}(\omega) = \begin{bmatrix} 1 & \hat{G}_{i,j}(\omega) \\ \hat{G}_{j,i}(\omega) & 1 \end{bmatrix}, \quad \text{where} \quad \Psi(t) = \begin{bmatrix} \delta(t) & G_{i,j}(t) \\ G_{j,i}(t) & \delta(t) \end{bmatrix} \quad (10)$$

is the impulse response matrix of the 2×2 system consisting of only nodes \mathbf{x}_i and \mathbf{x}_j , and $\hat{\Psi}(\omega)$ is Fourier transform (FT) of $\Psi[t]$. Also, $\hat{G}_{i,j}(\omega) \approx \pi^{-\frac{1}{2}} \exp\{-i\omega\mu_{i,j}\} \exp\{-\sigma_{i,j}^2\omega^2\} \otimes \mathcal{F}(\omega)$ is the FT of $G[t]_{i,j}$ while the FT of $\delta(t)$ is 1. The FT of approximation to $\hat{G}_{i,j}(\omega)$ is obtained by convolution of the FT of the multiplicative error term $\exp\{-\mathcal{O}(t^3)\}$ with the FT of $\exp\{-0.5(t - \mu^2\sigma^{-2})\} / \sqrt{2\pi\sigma^2}$.

However for stable VAR system, the transfer function between any pair of variables conditioned on the rest is [1]:

$$\hat{\Psi}(\omega) = \begin{bmatrix} \hat{A}_{i,i}(\omega) & \hat{A}_{i,j}(\omega) \\ \hat{A}_{j,i}(\omega) & \hat{A}_{j,j}(\omega) \end{bmatrix}^{-1}, \quad (11)$$

where $\hat{A}_{i,j}(\omega)$ is the FT of $A_{i,j}[t]$. Inverting the matrix in eqn. (10) and equating with terms of eqn. (11) gives

$$\hat{A}_{i,j}(\omega) = -\frac{\hat{G}_{i,j}(\omega)}{1 - \hat{G}(\omega)_{i,j}^* \hat{G}(\omega)_{i,j}}$$

Taking logs of the absolute value (squared) yields

$$\log(\hat{A}_{i,j}(\omega)^* \hat{A}_{i,j}(\omega)) = \log(\hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega)) - 2\log(1 - \hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega)).$$

However, since $|\hat{G}| \ll 1$ the square magnitude $\hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega) \ll 1$ which implies that

$$\log(1 - \hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega)) \approx 0 + \mathcal{O}(\hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega))$$

Substituting gives

$$\log(\hat{A}_{i,j}(\omega)^* \hat{A}_{i,j}(\omega)) \approx \log(\hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega)) + \mathcal{O}(\hat{G}_{i,j}(\omega)^* \hat{G}_{i,j}(\omega))$$

which implies that the FT of $|\hat{A}_{i,j}(\omega)| \approx |\hat{G}_{i,j}(\omega)|$ and therefore $A[t]_{i,j}$ can be approximated by a Gaussian function $\mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$. Moreover the approximation has a multiplicative error of order $\exp\{-\mathcal{O}(t^3)\}$. □

B.1 Proof for Proposition 2

Consider the dimensionless constant-coefficient convection-diffusion equation in 1-d:

$$\frac{\partial f}{\partial t} + \frac{\partial f}{\partial x} = \gamma \frac{\partial^2 f}{\partial t^2},$$

where f is the process, x is the spatial coordinate and $\gamma = 1/\text{Pe}$ is the inverse of the Péclet number of the system. Under infinite boundary conditions, the Green's function (*i.e.* impulse response) has the form

$$g(x, t) = \frac{1}{\sqrt{4\pi\gamma t}} \exp\left\{-\frac{1}{2} \frac{(x-t)^2}{2\gamma t}\right\}. \quad (12)$$

In order to derive an approximation, assume $\gamma = 1$ without loss of generality.

$$g(x, t) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{(x-t)^2}{2t} - \frac{1}{2} \log 2t\right\} \quad (13)$$

Writing $t = x + \tau$, a Taylor series expansion (TSE) of the term in the exponent gives:

$$\begin{aligned}
& \exp \left\{ -\frac{1}{2} \frac{\tau^2}{2(x+\tau)} - \frac{1}{2} \log 2(x+\tau) \right\} \\
&= \exp \left\{ -\frac{1}{2} \frac{\tau^2}{2} \left((2x)^{-1} - (2x)^{-2} \tau + (2x)^{-3} \tau^2 - (2x)^{-4} \tau^3 + \dots \right) \right. \\
&\quad \left. - \frac{1}{2} \left(\log 2x + (2x)^{-1} \tau - \frac{(2x)^{-2}}{2} \tau^2 + \frac{(2x)^{-3}}{3} \tau^3 + \dots \right) \right\} \\
&= \exp \left\{ -\frac{1}{2} \frac{\tau^2}{2x} - \frac{1}{2} \log 2x - \frac{1}{2} \left(\tau^3 \left((2x)^{-1} - (2x)^{-2} \right) - \tau^4 \left(\frac{(2x)^{-2}}{2} - (2x)^{-3} \right) + \dots \right) \right\} \\
&\approx \frac{1}{\sqrt{x}} \exp \left\{ -\frac{1}{2} \frac{\tau^2}{2x} \right\} \exp \left\{ -\frac{1}{2} \mathcal{O}(\tau^3) \right\}.
\end{aligned}$$

Therefore, for the far-field (*i.e.* $x \gg 0$) the approximation holds because:

$$\begin{aligned}
& \frac{\exp \left\{ -\frac{1}{2} \frac{\tau^2}{2x} \right\}}{\exp \left\{ -\frac{1}{2} \mathcal{O}(\tau^3) \right\}} \rightarrow 1 \quad \text{as} \quad \tau \rightarrow 0 \\
& \text{and} \quad \exp \left\{ -\frac{1}{2} \tau^3 \right\} - \exp \left\{ -\frac{1}{2} \frac{\tau^2}{2x} \right\} \rightarrow 0 \quad \text{as} \quad \tau \rightarrow \infty,
\end{aligned}$$

And for the near-field *i.e.* as $x \rightarrow 0$ the response function $g(x, t) \rightarrow \delta(t)$.

C Proof for Theorem 3

Proof. We start by introducing the following notation - let $\mathbf{y}_i \in \mathbf{x}_k$ imply that $Z_{i,k} = 1$. For two vertices \mathbf{y}_i and \mathbf{y}_j , let $\mathbf{y}_i \sim \mathbf{y}_j$ indicate that there is at least one shared latent component $\mathbf{y}_i \in \mathbf{x}_k$ and $\mathbf{y}_j \in \mathbf{x}_k$, *i.e.* $Z_{i,k} Z_{j,k} = 1$. Also, let $\hat{\mathbf{B}} = (\mathbf{Z}\mathbf{Z}^\top) \odot \mathbf{B}$, where \odot is the Hadamard product. Therefore, $\hat{\mathbf{B}}_{i,j} = 0 \Rightarrow \exists k$ s.t. $\mathbf{y}_i \sim \mathbf{y}_j$. Without loss of generality, we will prove this assertion for the case below, that is:

Proposition 3. For a given \mathbf{Z} , the least-squares (LS) local optimum \mathbf{A}^* and \mathbf{x}^* to

$$\mathbf{x}[t] = \sum_{p=1}^P \mathbf{A}[p] \mathbf{x}[t-p] + \mathbf{u}[t] \quad (14)$$

$$\mathbf{y}[t] = \hat{\mathbf{B}} \mathbf{y}[t-1] + \mathbf{Z} \mathbf{x}[t] + \mathbf{v}[t], \quad (15)$$

is also a local optimum for

$$\mathbf{x}[t] = \sum_{p=1}^P \mathbf{A}[p] \mathbf{x}[t-p] + \mathbf{u}[t] \quad (16)$$

$$\mathbf{y}[t] = \mathbf{C} \mathbf{y}[t-1] + \mathbf{Z} \mathbf{x}[t] + \mathbf{v}[t], \quad (17)$$

for some diagonal matrix \mathbf{C} .

The more general case of Theorem 3 is a straightforward extension of this proof.

First, we observe that as eqn. (14) is decoupled from \mathbf{B} , the local minimum $\mathbf{A}^*[p]$ conditioned on \mathbf{x}^* is de-coupled from \mathbf{B} . Therefore, we only need to show that if

$$\mathbf{x}^+ = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x})$$

where:

$$f(\mathbf{x}) = \sum_t \left[\mathbf{y}[t] - \hat{\mathbf{B}}\mathbf{y}[t-1] - \mathbf{Z}\mathbf{x}[t] \right]^\top \left[\mathbf{y}[t] - \mathbf{B}\mathbf{y}[t-1] - \mathbf{Z}\mathbf{x}[t] \right]$$

then \mathbf{x}^+ also a LS solution to eqn. (17).

Now,

$$\frac{1}{2} \nabla_{\mathbf{x}[t]} f(\mathbf{x}) = \left[\mathbf{y}[t] - \mathbf{B}\mathbf{y}[t-1] - \mathbf{Z}\mathbf{x}[t] \right]^\top \mathbf{Z} = \mathbf{y}[t]^\top \mathbf{Z} - \mathbf{y}[t-1]^\top \hat{\mathbf{B}}^\top \mathbf{Z} - \mathbf{x}[t]^\top \mathbf{Z}^\top \mathbf{Z}$$

Defining $\eta = \hat{\mathbf{B}}\mathbf{y}[t-1]$, we get $\eta_i = [\hat{\mathbf{B}}\mathbf{y}[t-1]]_i = \sum_j \mathbf{y}_j[t-1] \hat{B}_{i,j}$, that is $\eta_i = \sum_j B_{i,j} \mathbf{y}_j[t-1]$ for all \mathbf{y}_j such that $\mathbf{y}_i \sim \mathbf{y}_j$.

Moreover, $[\eta^\top \mathbf{Z}]_k = \sum_i \eta_i Z_{i,k}$ that is $[\eta^\top \mathbf{Z}]_k = \sum_i \eta_i$ for all $\mathbf{y}_i \in \mathbf{x}_k$. Therefore, chaining these two we get $[\mathbf{Z}^\top \eta]_k = \sum_i \sum_j B_{i,j} \mathbf{y}_j[t-1]$ for all i s.t. $\mathbf{y}_i \in \mathbf{x}_k$ and j s.t. $\mathbf{y}_i \sim \mathbf{y}_j$, which is the same as $[\mathbf{Z}^\top \eta]_k = \sum_{\mathbf{y}_i \in \mathbf{x}_k} \sum_{\mathbf{y}_j \in \mathbf{x}_k} B_{i,j} \mathbf{y}_j[t-1]$.

Therefore

$$\frac{1}{2} [\nabla_{\mathbf{x}[t]} f(\mathbf{x})]_k = [\mathbf{y}[t]^\top \mathbf{Z}]_k - [\mathbf{Z}^\top \eta]_k = \sum_{\mathbf{y}_i \in \mathbf{x}_k} \sum_{\mathbf{y}_j \in \mathbf{x}_k} B_{i,j} \mathbf{y}_j[t-1] - [\mathbf{x}[t]^\top \mathbf{Z}^\top \mathbf{Z}]_k \quad (18)$$

This gradient is equivalent to $\frac{1}{2} \nabla_{\mathbf{x}[t]} g(\mathbf{x})$ where $g(\mathbf{x})$ is the LS objective of eqn. (17).

$$\frac{1}{2} \nabla_{\mathbf{x}[t]} g(\mathbf{x}) = \mathbf{y}[t]^\top \mathbf{Z} - \mathbf{y}[t-1]^\top \mathbf{C} \mathbf{Z} - \mathbf{x}[t]^\top \mathbf{Z}^\top \mathbf{Z}$$

where $[\mathbf{y}[t-1]^\top \mathbf{C} \mathbf{Z}]_k = \sum_{j \in \mathbf{x}_k} C_j \mathbf{y}_j[t-1]$, for $C_j = \sum_{\mathbf{y}_i \in \mathbf{x}_k} B_{i,j}$.

Note that the Hessian of the objective function is independent of B and C. Therefore, the two quadratic problems differ only in a constant independent of \mathbf{x} and therefore have the same optimal solutions.

□

D Optimization with respect to \mathbf{x}

The state-space model

$$\mathbf{x}[t] = \sum_{p=1}^P \mathbf{A}[p] \mathbf{x}[t-p] + \mathbf{u}[t] \quad \text{and} \quad \mathbf{y}[t] = \sum_{q=1}^Q \mathbf{C}[q] \mathbf{y}[t-q] + \mathbf{Z}\mathbf{x}[t] + \mathbf{v}[t], \quad (19)$$

can be re-written as the followed augmented state space form

$$\zeta[t+1] = \mathbf{A}_{\text{aug}} \zeta'[t] + \nu[t] \quad \text{and} \quad \vartheta[t] = \mathbf{Z}_{\text{aug}} \zeta'[t] + \mathbf{v}[t],$$

where $\zeta[t] = (\mathbf{x}[t] \dots \mathbf{x}[t-P])^\top$ is the augmented state, $\vartheta[t] = (\mathbf{y}[t] - \sum_{q=1}^Q \mathbf{C}^{(n)}[q] \mathbf{y}[t-q])^\top$ is the augmented observation, $\nu[t] = (\mathbf{u}[t], 0 \dots 0)^\top$ is the augmented state innovations the matrices \mathbf{A}_{aug} is a $PK \times PK$ matrix with $\mathbf{A}^{(n)}[1] \dots \mathbf{A}^{(n)}[P]$ on the first row and \mathbf{Z}_{aug} is a $N \times KP$ matrix constructed from $\mathbf{A}^{(n)}$ and $\mathbf{Z}^{(n)}$. Namely:

$$\mathbf{A}_{\text{aug}} = \begin{bmatrix} \mathbf{A}[1] & \mathbf{A}[2] & \dots & \mathbf{A}[P] \\ \mathbf{I} & 0 & \dots & 0 \\ 0 & \mathbf{I} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad \mathbf{Z}_{\text{aug}} = [\mathbf{Z} \quad 0 \quad \dots \quad 0]$$

Rewriting the summation of time in vector-matrix format for notational clarity, the optimization problem is then:

$$\zeta^{(n+1)} = \underset{\zeta}{\operatorname{argmin}} \|\vartheta - \Delta_{\text{aug}} \zeta\|^2 + \lambda_0 \|\Lambda_{\text{aug}} \zeta\|^2$$

where Δ and Λ are the $T \times T$ block-diagonal matrices:

$$\Delta = \begin{bmatrix} Z_{\text{aug}} & 0 & \dots & 0 \\ 0 & Z_{\text{aug}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & Z_{\text{aug}} \end{bmatrix} \quad \Lambda_{\text{aug}} = \begin{bmatrix} -\Lambda_{\text{aug}} & I & 0 & \dots & 0 \\ 0 & -\Lambda_{\text{aug}} & I & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & -\Lambda_{\text{aug}} \end{bmatrix}$$

This is a $T \times T$ tri-diagonal Toeplitz system with blocks of size KP and which can be solved using specialized solvers that have running time of $\mathcal{O}(T \times KP^3)$ [2]. In case the system is ill-conditioned or rank-deficient, a small perturbation may be added along the diagonal of the matrix. However in practice, we have observed the system to be well conditioned, especially when close to a solution point.

E Proximal operators

E.1 Operator for Z

Defining $Z_{i,\cdot} = \{Z_{i,1} \dots Z_{i,K}\}$, the proximal operator $\mathbf{prox}_h(Z^{(n)})$ with respect to Z , and writing the unit-ball constraint $\|Z_{i,\cdot}\|_2 \leq 1$ as an indicator function $\mathbb{I}_1(\|Z_{i,\cdot}\|_2)$, we get:

$$\begin{aligned} \mathbf{prox}_g(\hat{Z}) &= \underset{Z}{\operatorname{argmin}} \lambda_3 \|Z\|_1 + \sum_{i=1}^N \mathbb{I}_1(\|Z_{i,\cdot}\|_2) + \frac{1}{2} \|Z - \hat{Z}\|^2 \\ &= \lambda_3 \sum_{i=1}^N \|Z_{i,\cdot}\|_1 + \sum_{i=1}^N \mathbb{I}_1(\|Z_{i,\cdot}\|_2) + \frac{1}{2} \sum_{i=1}^N \|Z_{i,\cdot} - \hat{Z}_{i,\cdot}\|^2 \end{aligned}$$

As the problem is decomposable into a sum of problems over $Z[1, \cdot] \dots Z[N, \cdot]$, consider the individual problem of the form:

$$\begin{aligned} &\underset{\mathbf{z}}{\operatorname{argmin}} \lambda_3 \|\mathbf{z}\|_1 + \frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|^2 \\ &\text{subject to} \quad \|\mathbf{z}\|_2 \leq 1 \end{aligned}$$

where $\mathbf{z} = \{\mathbf{z}_1 \dots \mathbf{z}_K\}$ represents any $Z_{i,\cdot}$.

Introducing Lagrange multipliers $\mu \in \mathbb{R}^+$ and $\eta \in \mathbb{R}^{+K}$, the dual is:

$$\mathcal{L}(\mathbf{z}, \mu, \eta) = \lambda_3 \|\mathbf{z}\|_1 + \frac{1}{2} \|\mathbf{z} - \hat{\mathbf{z}}\|^2 + \frac{\mu}{2} (\|\mathbf{z}\|_2^2 - 1)$$

Dual feasibility implies:

$$0 = \partial \mathcal{L}(\mathbf{z}, \mu, \eta)_{\mathbf{z}_k} = \lambda_3 \partial_k \|\mathbf{z}\|_1 + \mathbf{z}_k - \hat{\mathbf{z}}_k + \mu \mathbf{z}_k$$

therefore in vector format

$$\mathbf{z}(1 + \mu) = \hat{\mathbf{z}} - \lambda_3 \partial \|\mathbf{z}\|_1$$

where the sub-gradient

$$\partial_k \|\mathbf{z}\|_1 = \begin{cases} 1 & \text{if } \mathbf{z}_k > 0 \\ -1 & \text{if } \mathbf{z}_k < 0 \\ (-1, +1) & \text{if } \mathbf{z}_k = 0 \end{cases}$$

Now if there is at-least one k such that $|\hat{\mathbf{z}}_k| > \lambda$, then the following argument applies;

If $\mathbf{z}_k^* > 0$ then $\partial \|\mathbf{z}\|_1 = 1$. Therefore, $(1 + \mu)\mathbf{z}_k^* = (\mathbf{z}' - \lambda) > 0$. This implies that if $\mathbf{z}' > \lambda$ then $(1 + \mu)\mathbf{z}_k^* = (\hat{\mathbf{z}} - \lambda)$. Similarly if $\hat{\mathbf{z}} > -\lambda$ then $(1 + \mu)\mathbf{z}_k^* = (\hat{\mathbf{z}} + \lambda)$, and if $-\lambda \leq \hat{\mathbf{z}} \leq \lambda$ then $\mathbf{z}_k^* = 0$.

Defining T_λ as the element-wise shrinkage operator, we get $(1 + \mu)\mathbf{z}^* = T_\lambda(\mathbf{z}')$. Moreover, setting $\mu = \min(\|T_\lambda(\mathbf{z}')\|_2 - 1, 0)$ satisfies the complementary slackness condition:

$$\begin{aligned} \mu^* > 0 & \quad \text{if} \quad \|\mathbf{z}\|_2^2 = 1 \\ \mu^* = 0 & \quad \text{if} \quad \|\mathbf{z}\|_2^2 < 1 \end{aligned}$$

However if $|\hat{\mathbf{z}}_k| \leq \lambda$, then $\mathbf{z}^* = 0$ is the only feasible solution.

Therefore, the final solution is:

$$\mathbf{z}^* = \max\left(1, \frac{1}{\|T_\lambda(\hat{\mathbf{z}})\|_2}\right) T_\lambda(\hat{\mathbf{z}})$$

E.2 Operator for A

The proximal operator for A is

$$\mathbf{prox}_g(\hat{\mathbf{A}}) = \underset{\mathbf{A}}{\operatorname{argmin}} \sum_{i,j=1}^K \left[\lambda_1 \|\mathbf{D}(\gamma_{i,j})\mathbf{A}_{i,j}\|_2 + \frac{1}{2} \|\mathbf{A}_{i,j} - \widehat{\mathbf{A}}_{i,j}\|^2 \right],$$

subject to $0 \leq \sum_p \mathbf{A}_{i,j}[p] \leq 1$ for all $i, j = 1 \dots K$, where $\mathbf{D}(\gamma_{i,j})$ represents the operator $\hat{\partial}_t + \gamma_{i,j}(p - \mu_{i,j})$ and is a $P \times P$ full-rank matrix. Since this can be split across all $\mathbf{A}_{i,j}$, consider one problem of the form:

$$\mathbf{prox}_g(\hat{\mathbf{a}}) = \underset{\mathbf{a}}{\operatorname{argmin}} \lambda_1 \|\mathbf{D}\mathbf{a}\|_2 + \frac{1}{2} \|\mathbf{a} - \hat{\mathbf{a}}\|^2, \quad (20)$$

subject to $0 \leq \mathbf{1}^\top \mathbf{a} \leq 1$ where \mathbf{a} represents any $\mathbf{A}_{i,j}$, and \mathbf{D} represents the corresponding $\mathbf{D}(\gamma_{i,j})$.

The Lagrangian of the problem is:

$$\mathcal{L}_{\mathbf{prox}_g}(\mathbf{a}, \eta_1, \eta_2) = \lambda_1 \|\mathbf{D}\mathbf{a}\|_2 + \frac{1}{2} \|\mathbf{a} - \hat{\mathbf{a}}\|^2 - \eta_1 \mathbf{1}^\top \mathbf{a} + \eta_2 (\mathbf{1}^\top \mathbf{a} - 1) \quad (21)$$

where $\eta_1 > 0$ and $\eta_2 > 0$ are Lagrange multipliers for $-\mathbf{1}^\top \mathbf{a} \leq 0$ and $\mathbf{1}^\top \mathbf{a} \leq 1$

As \mathbf{D} is full rank, define $\mathbf{b} = \mathbf{D}\mathbf{a}$ and $\hat{\mathbf{b}} = \mathbf{D}^{-1}\hat{\mathbf{a}}$. Therefore, eqn. (22) can be rewritten as :

$$\mathbf{prox}_g(\hat{\mathbf{a}}) = \mathbf{D} \left[\underset{\mathbf{b}}{\operatorname{argmin}} \lambda_1 \|\mathbf{b}\|_2 + \frac{1}{2} \|(\mathbf{b} - \hat{\mathbf{b}})\|^2 \right],$$

subject to $0 \leq \mathbf{1}^\top \mathbf{D}^{-1}\mathbf{b} \leq 1$.

The dual feasible solution $\mathbf{b}^*(\eta_1, \eta_2)$ of the corresponding Lagrangian satisfies the equation:

$$\mathbf{b}^*(\eta_1, \eta_2) = \hat{\mathbf{b}} - \lambda_1 \partial_{\mathbf{b}} \|\mathbf{b}^*\|_2 + (\eta_1 - \eta_2) \mathbf{D}^{-1} \mathbf{1},$$

where $\partial_{\mathbf{b}} \|\mathbf{D}\mathbf{b}^*\|_2$, the sub-gradient of $\sqrt{\sum_{p'=1}^P (\sum_{p=1}^P \mathbf{b}[p])^2}$ is given by:

$$\partial_{\mathbf{b}[p]} \|\mathbf{b}^*\|_2 = \begin{cases} \frac{\mathbf{b}^*[p]}{\|\mathbf{b}^*\|_2} & \text{if } \|\mathbf{b}^*\|_2 > 0 \\ (-1, +1) & \text{if } \|\mathbf{b}^*\|_2 = 0 \end{cases}$$

Therefore, given values of η_1 and η_2 and defining $\beta = (\eta_1 - \eta_2) \mathbf{D}^{-1} \mathbf{1}$, we get

$$\mathbf{b}^*[p](\eta_1, \eta_2) = \begin{cases} \hat{\mathbf{b}}[p] + \beta[p] - \lambda \frac{\hat{\mathbf{b}}[p] + \beta[p]}{\|\hat{\mathbf{b}} + \beta\|_2}, & \text{if } \|\hat{\mathbf{b}} + \beta\|_2 > \lambda \\ 0 & \text{otherwise} \end{cases}$$

which gives:

$$\mathbf{a}^*(\eta_1, \eta_2) = \begin{cases} \hat{\mathbf{a}} + (\eta_1 - \eta_2)\mathbf{1} - \lambda \frac{\hat{\mathbf{a}} + (\eta_1 - \eta_2)\mathbf{1}}{\|\mathbf{D}^{-1}\hat{\mathbf{a}} + (\eta_1 - \eta_2)\mathbf{1}\|_2} & \text{if } \|\mathbf{D}^{-1}\hat{\mathbf{a}} + (\eta_1 - \eta_2)\mathbf{1}\|_2 > \lambda \\ 0 & \text{otherwise} \end{cases} \quad (22)$$

And the gradient of eqn. (21) with respect the dual variables at the dual-feasible solution is :

$$\begin{aligned} \nabla_{\eta_1} \mathcal{L} &= -\mathbf{1}^\top \mathbf{a} \\ \text{and } \nabla_{\eta_2} \mathcal{L} &= \mathbf{1}^\top \mathbf{a} \end{aligned}$$

Therefore, the dual-ascent step is

$$\begin{aligned} \eta_1^{(k+1)} &= \eta_1^{(k)} - \alpha^{(k)} \min\{\mathbf{1}^\top \mathbf{a}, 0\} \\ \eta_2^{(k+1)} &= \eta_2^{(k)} + \alpha^{(k)} \max\{\mathbf{1}^\top \mathbf{a} - 1, 0\} \end{aligned}$$

Therefore using the fact that violation of constraints for η_1 and η_2 are mutually exclusive, we combine η_1, η_2 into a single dual variable $\eta \in \mathbb{R}$, giving the the dual feasible solution as:

$$\mathbf{a}^{(n+1)} = \begin{cases} \hat{\mathbf{a}} + \eta^{(n)}\mathbf{1} - \lambda \frac{\hat{\mathbf{a}} + \eta^{(n)}\mathbf{1}}{\|\mathbf{D}^{-1}\hat{\mathbf{a}} + \eta^{(n)}\mathbf{1}\|_2} & \text{if } \|\mathbf{D}^{-1}\hat{\mathbf{a}} + \eta^{(n)}\mathbf{1}\|_2 > \lambda \\ 0 & \text{otherwise} \end{cases} \quad (23)$$

and the dual-ascent step is:

$$\eta^{(n+1)} = \begin{cases} \eta^{(n)} - \alpha^{(n)} \mathbf{1}^\top \mathbf{a}^{(n)} & \text{if } \mathbf{1}^\top \mathbf{a}^{(n)} < 0 \\ \eta^{(n)} + \alpha^{(n)} (\mathbf{1}^\top \mathbf{a}^{(n)} - 1) & \text{if } \mathbf{1}^\top \mathbf{a}^{(n)} > 1 \end{cases} \quad (24)$$

F Dependence on Initialization

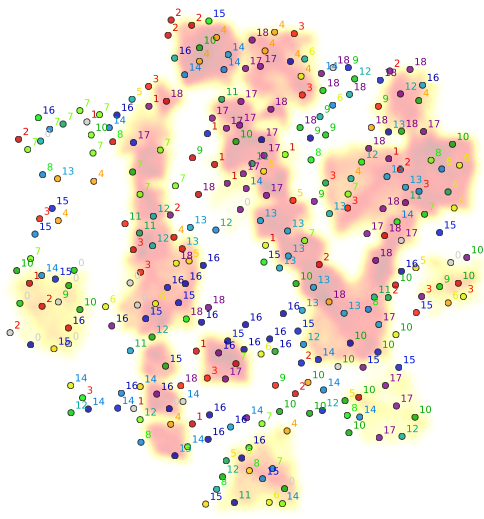
In order to demonstrate the robustness of the solution to initialization, here we show the result of estimation procedure for different initializations of the K -means step. As mentioned earlier, the initial positions of the K -means centroids are uniformly distributed at random. It can be observed that regardless of the K -means solution the algorithm converges to highly consistent global graphical structures. Note that for each solution, labels values have been selected to maximize overlap across results to improve comparison.

G Cross-Validation

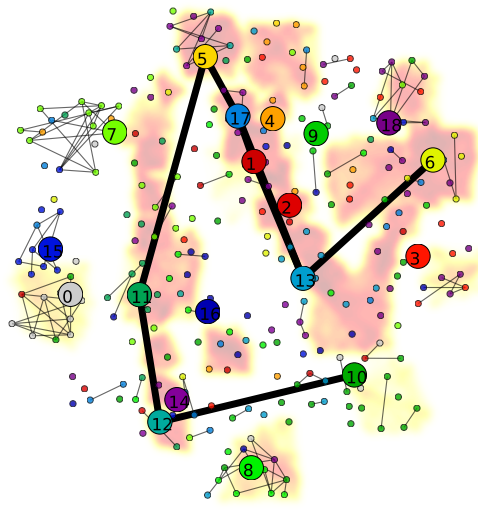
A 10-fold block cross validation approach is used for model selection and to assess the performance of the various models. Here the time-series $\mathbf{y}[t]; 1 = 1 \dots T$ is divided into 10 contiguous blocks of length $T/10$ each and the model parameters estimated using 9 blocks. The data from the remaining block (indexed as $t = 1 \dots T'$) are the fitted to the model using least squares to get $\hat{\mathbf{y}}$. For example for the standard order- P VAR model with parameters $\mathbf{A}[1] \dots \mathbf{A}[P]$, this would be

$$\begin{aligned} &\min_{\hat{\mathbf{y}}} \sum_t \|\hat{\mathbf{y}}[t] - \mathbf{y}[t]\|_2 \\ \text{such that } &\hat{\mathbf{y}}[t] = \sum_{p=1}^P \mathbf{A}[p] \hat{\mathbf{y}}[t-p] \end{aligned}$$

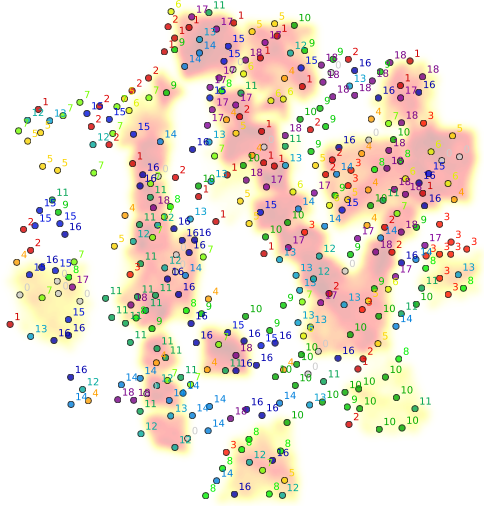
Defining $\eta[t]; t = 1 \dots T$ as the Lagrange multiplier for each constraint and $\rho > 0$ as the augmented Lagrangian multiplier term, the solution for this is computed using a dual-ascent procedure, where



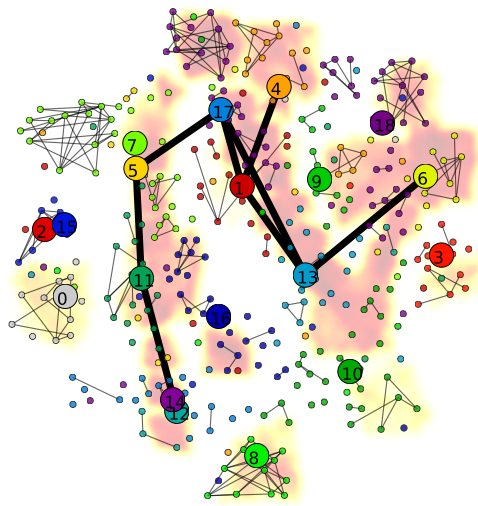
(a) K-Means initialization



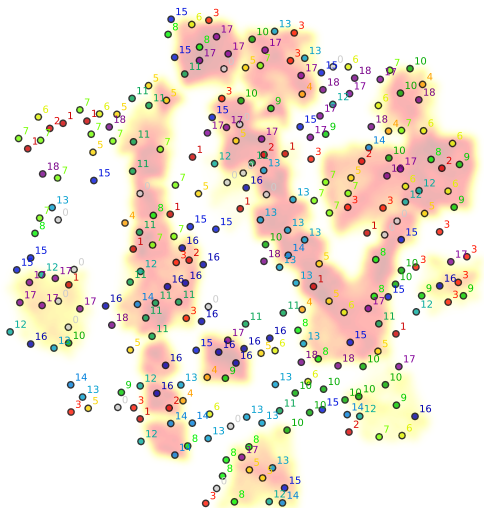
(b) Estimated graphical structure



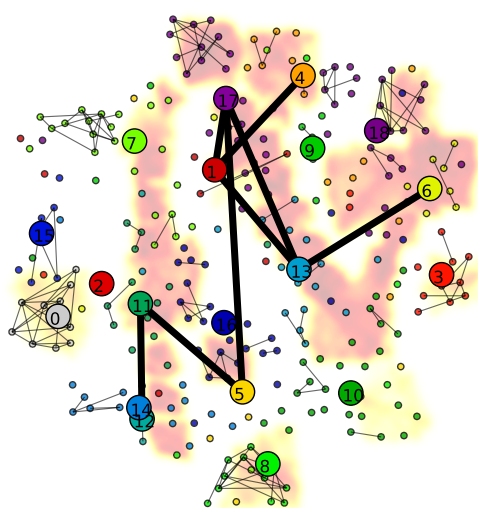
(c) K-Means initialization



(d) Estimated graphical structure



(e) K-Means initialization



(f) Estimated graphical structure

Figure 1: The left column shows the node assignments after one K -means initialization step. The right column shows the estimated graphical structure for the corresponding initialization. The nodes in the latent global graph are positioned at the centroids of the corresponding local graphs. In all figures, colors and labels both indicate cluster assignments of the sites. Labels for the local nodes are omitted for clarity.

each iteration involves the following two steps:

$$\begin{aligned} \text{[Primal update]: } \hat{\mathbf{y}}[t]^{(k+1)} &= \mathbf{y}[t] - \frac{1}{2} \left(\eta^{(k)}[t] - \sum_{p=1}^P \mathbf{A}[p]^\top \eta^{(k)}[t+p] \right) \\ \text{[Dual update]: } \eta[t]^{(k+1)} &= \eta[t]^{(k)} + \rho \left(\hat{\mathbf{y}}^{(k+1)}[t] - \sum_{p=1}^P \mathbf{A}[p] \hat{\mathbf{y}}^{(k+1)}[t-p] \right) \end{aligned}$$

Estimating the best least-squares fit for the hierarchical model, given parameters \mathbf{A} , \mathbf{B} and \mathbf{Z} and hyper-parameter λ_0 , involves solving

$$\begin{aligned} \min_{\hat{\mathbf{y}}, \mathbf{x}} \sum_t \|\hat{\mathbf{y}}[t] - \mathbf{y}[t]\|_2 + \left\| \mathbf{x}[t] - \lambda_0 \sum_{p=1}^P \mathbf{A}[p] \mathbf{x}[t-p] \right\|_2 \\ \text{such that } \hat{\mathbf{y}}[t] = \sum_{q=1}^Q \mathbf{B}[q] \hat{\mathbf{y}}[t-p] + \mathbf{Z} \mathbf{x}[t] \end{aligned}$$

which can be solved using the following dual-ascent scheme:

$$\begin{aligned} \text{[Primal update]: } \hat{\mathbf{y}}[t]^{(k+1)} &= \mathbf{y}[t] - \frac{1}{2} \left(\eta^{(k)}[t] - \sum_{q=1}^Q \mathbf{B}[q]^\top \eta^{(k)}[t+p] \right) \\ \text{[Dual update]: } \eta[t]^{(k+1)} &= \eta[t]^{(k)} + \rho \left(\hat{\mathbf{y}}^{(k+1)}[t] - \sum_{q=1}^Q \mathbf{B}[q] \hat{\mathbf{y}}^{(k+1)}[t-p] - \mathbf{Z} \mathbf{x}^{(k+1)}[t] \right) \end{aligned}$$

while $\mathbf{x}^{(k+1)}$ is estimated using the method described in Supplemental Appendix [D](#).

The relative error is then defined as

$$\sqrt{\frac{\sum_t \|\mathbf{y}[t] - \hat{\mathbf{y}}[t]\|}{\sum_t \|\mathbf{y}[t]\|}}$$

References

- [1] Brockwell, P.J., Davis, R.A.: Time Series: Theory and Methods. Springer Series in Statistics, 2 edn. (1991) [2](#), [4](#)
- [2] Minchev, B.V.: Some algorithms for solving special tridiagonal block toeplitz linear systems. Journal of Computational and Applied Mathematics 156, 179–200 (2003) [7](#)
- [3] Schutter, B.D.: Minimal state-space realization in linear system theory: An overview. Tech. report bds:99-07, Delft University of Technology (2000) [2](#)