

A An Illustrative Example of Mori-Zwanzig Formalism

The following problem is widely used as an introductory example of the Mori-Zwanzig (MZ) formalism [21, 51] for model order reduction: let $\mathbf{z} = [\mathbf{x}; \mathbf{y}] \in \mathbb{R}^N$ where $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^{N-n}$, and $n \ll N$. Consider the system of linear differential equations

$$\begin{cases} \mathbf{x}' = \mathbf{A}_{11}\mathbf{x} + \mathbf{A}_{12}\mathbf{y}, \\ \mathbf{y}' = \mathbf{A}_{21}\mathbf{x} + \mathbf{A}_{22}\mathbf{y}, \end{cases} \quad (22)$$

with initial values $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{y}(0) = \mathbf{y}_0$. Suppose we are interested in the time evolution of $\mathbf{x}(t)$, which depends on the joint effect of \mathbf{x} and \mathbf{y} . However the computation of the complete system (22) is expensive and can be prohibitive for large N . The question is whether we can derive a reduced system only involving \mathbf{x} from (22). To this end, we assume \mathbf{x} is given, and solve for \mathbf{y} from the \mathbf{y} -equation of (22) to obtain

$$\mathbf{y}(t) = e^{\mathbf{A}_{22}t}\mathbf{y}_0 + \int_0^t e^{\mathbf{A}_{22}(t-s)}\mathbf{A}_{22}\mathbf{x}(s) ds. \quad (23)$$

Then we plug this back into the \mathbf{x} -equation of (22) and obtain

$$\mathbf{x}'(t) = \mathbf{A}_{11}\mathbf{x}(t) + \mathbf{A}_{12} \int_0^t e^{\mathbf{A}_{22}(t-s)}\mathbf{A}_{22}\mathbf{x}(s) ds + \mathbf{A}_{12}e^{\mathbf{A}_{22}t}\mathbf{y}_0, \quad (24)$$

which neglects the dependence on $\mathbf{y}(t)$ except for the initial value \mathbf{y}_0 .

As shown in the example above, MZ formalism aims at reducing a high-dimensional system of \mathbf{z} into a low-dimensional system of \mathbf{x} (resolved variable) while maintaining the effect of \mathbf{y} (unresolved variable). This is particularly useful if an exact solution of \mathbf{z} is unnecessary to understand the dynamics of \mathbf{x} . Specialized derivations and subsequent approximation techniques can be implemented to obtain highly efficient numerical solutions for nonlinear systems.

B Proofs

B.1 Proof of Theorem 1

Proof. Let $\lambda_i^*(t)$ be the conditional intensity of node i at time t , i.e., $\mathbb{E}[dX_i(t)|\mathcal{H}(t)] = \lambda_i^*(t) dt$. In the standard diffusion model, the conditional intensity $\lambda_i^*(t)$ of a healthy node i (i.e., $X_i(t) = 0$) is determined by the total infection rate of its infected neighbors j (i.e., $X_j(t) = 1$). That is,

$$\lambda_i^*(t) = \sum_j \alpha_{ji}X_j(t)(1 - X_i(t)). \quad (25)$$

By taking expectation $\mathbb{E}_{\mathcal{H}(t)}[\cdot]$ on both sides of (25), we obtain

$$\begin{aligned} \lambda_i(t) &:= \mathbb{E}_{\mathcal{H}(t)}[\lambda_i^*(t)] = \mathbb{E}_{\mathcal{H}(t)} \left[\alpha_{ji}X_j(t)(1 - X_i(t)) | \mathcal{H}(t) \right] \\ &= \sum_j \alpha_{ji}(x_j - x_{ij}) = \sum_j \alpha_{ji}(x_j - y_{ij} - e_{ij}). \end{aligned} \quad (26)$$

On the other hand, there is

$$\lambda_i(t) dt = \mathbb{E}_{\mathcal{H}(t)}[dX_i^*(t)] = \mathbb{E}_{\mathcal{H}(t)}[dX_i(t)|\mathcal{H}(t)] = d\mathbb{E}_{\mathcal{H}(t)}[X_i(t)|\mathcal{H}(t)] = dx_i. \quad (27)$$

Combining (26) and (27) yields

$$x'_i = \frac{dx_i(t)}{dt} = \sum_j \alpha_{ji}(x_j - y_{ij} - e_{ij}) = (\mathbf{A}\mathbf{x})_i - (\text{diag}(\mathbf{x})\mathbf{A}\mathbf{x})_i - \sum_j \alpha_{ji}e_{ij}$$

for every $i \in [n]$, which verifies the \mathbf{x} part of (6). Similarly, we can obtain

$$x'_I = \sum_{i \in I} \sum_{j \notin I} \alpha_{ji}(x_I - x_{I \cup \{j\}}) = \sum_{i \in I} \sum_{j \notin I} \alpha_{ji}(y_I + e_I - y_{I \cup \{j\}} - e_{I \cup \{j\}}). \quad (28)$$

Moreover, by taking derivative on both sides of $x_I(t) = y_I(t) + e_I(t)$, we obtain

$$x'_I = \sum_{i \in I} y_{I \setminus \{i\}} x'_i + e'_I = \sum_{i \in I} y_{I \setminus \{i\}} \sum_{j \neq i} \alpha_{ji} (x_j - x_i x_j - e_{ij}) + e'_I. \quad (29)$$

Combining (28) and (29) yields the e part of (6).

It is clear that $\mathbf{x}_0 = \chi_S$. For every I , at time $t = 0$, there is $x_I(0) = \prod_{i \in I} X_i(0) = 1$ if $I \subset S$ and 0 otherwise; and the same for $y_I(0)$. Hence $e_I(0) = x_I(0) - y_I(0) = 0$ for all I . Hence $\mathbf{z}_0 = [\mathbf{x}_0; \mathbf{e}_0] = [\chi_S; \mathbf{0}]$, which verifies the initial condition of (6). \square

B.2 Proof of Theorem 2

Proof. Consider the system (6) over a finite time horizon $[0, T]$, which evolves on a smooth manifold $\Gamma \subset \mathbb{R}^N$. For any real-valued phase (observable) space function $g : \Gamma \rightarrow \mathbb{R}$, the nonlinear system (6) is equivalent to the linear partial differential equation, known as the Liouville equation:

$$\begin{cases} \partial_t u(t, \mathbf{z}) = \mathcal{L}[u](t, \mathbf{z}), \\ u(0, \mathbf{z}) = g(\mathbf{z}), \end{cases} \quad (30)$$

where the Liouville operator $\mathcal{L}[u] := \bar{\mathbf{f}}(\mathbf{z}) \cdot \nabla_{\mathbf{z}} u$. The equivalency is in the sense that the solution of (30) satisfies $u(t, \mathbf{z}_0) = g(\mathbf{z}(t; \mathbf{z}_0))$, where $\mathbf{z}(t; \mathbf{z}_0)$ is the solution to (6) with initial value \mathbf{z}_0 .

Denote $e^{t\mathcal{L}}$ the Koopman operator associated with \mathcal{L} such that $e^{t\mathcal{L}}g(\mathbf{z}_0) = g(\mathbf{z}(t))$ where $\mathbf{z}(t)$ is the solution of (6). Then $e^{t\mathcal{L}}$ satisfies the semi-group property, i.e.,

$$e^{t\mathcal{L}}g(\mathbf{z}) = g(e^{t\mathcal{L}}\mathbf{z}) \quad (31)$$

for all g . On the right hand side of (31), \mathbf{z} can be interpreted as $\mathbf{z} = \iota(\mathbf{z}) = [\iota_1(\mathbf{z}), \dots, \iota_N(\mathbf{z})]$ where $\iota_j(\mathbf{z}) = z_j$ for all j .

Now consider the projection operator \mathcal{P} as the truncation such that $\mathcal{P}g(\mathbf{z}) = \mathcal{P}g(\mathbf{x}, \mathbf{e}) = g(\mathbf{x}, 0)$ for any $\mathbf{z} = (\mathbf{x}, \mathbf{e})$, and its orthogonal complement as $\mathcal{Q} = I - \mathcal{P}$ where I is the identity operator. Note that $\mathbf{z}'(t) = \frac{d\mathbf{z}(t)}{dt} = \frac{\partial}{\partial t} e^{t\mathcal{L}} \mathbf{z}_0$, and $\bar{\mathbf{f}}(\mathbf{z}(t)) = e^{t\mathcal{L}} \bar{\mathbf{f}}(\mathbf{z}_0) = e^{t\mathcal{L}} \mathcal{L} \mathbf{z}_0$ since $\mathcal{L}\iota_j(\mathbf{z}) = \mathbf{f}_j(\mathbf{z})$ for all \mathbf{z} and j . Therefore (6) implies that

$$\frac{\partial}{\partial t} e^{t\mathcal{L}} \mathbf{z}_0 = e^{t\mathcal{L}} \mathcal{L} \mathbf{z}_0 = e^{t\mathcal{L}} \mathcal{P} \mathcal{L} \mathbf{z}_0 + e^{t\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0. \quad (32)$$

Note that the first term on the right hand side of (32) is

$$e^{t\mathcal{L}} \mathcal{P} \mathcal{L} \mathbf{z}_0 = \mathcal{P} \mathcal{L} e^{t\mathcal{L}} \mathbf{z}_0 = \mathcal{P} \mathcal{L} \mathbf{z}(t). \quad (33)$$

For the second term in (32), we recall that the well-known Dyson's identity for the Koopman operator \mathcal{L} is given by

$$e^{t\mathcal{L}} = e^{t\mathcal{Q}\mathcal{L}} + \int_0^t e^{s\mathcal{L}} \mathcal{P} \mathcal{L} e^{(t-s)\mathcal{Q}\mathcal{L}} ds. \quad (34)$$

Applying (34) to $\mathcal{Q}\mathcal{L}\mathbf{z}_0$ yields

$$\begin{aligned} e^{t\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0 &= e^{t\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0 + \int_0^t e^{s\mathcal{L}} \mathcal{P} \mathcal{L} e^{(t-s)\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0 ds \\ &= e^{t\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0 + \int_0^t \mathcal{P} \mathcal{L} e^{(t-s)\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} e^{s\mathcal{L}} \mathbf{z}_0 ds \\ &= e^{t\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0 + \int_0^t \mathcal{P} \mathcal{L} e^{(t-s)\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}(s) ds. \end{aligned} \quad (35)$$

Substituting (33) and (35) into (32), we obtain

$$\frac{\partial}{\partial t} e^{t\mathcal{L}} \mathbf{z}_0 = \mathcal{P} \mathcal{L} \mathbf{z}(t) + e^{t\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}_0 + \int_0^t \mathcal{P} \mathcal{L} e^{(t-s)\mathcal{Q}\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}(s) ds, \quad (36)$$

where we used the fact that $e^{t\mathcal{L}} \mathcal{P} \mathcal{L} \mathbf{z}_0 = \mathcal{P} \mathcal{L} e^{t\mathcal{L}} \mathbf{z}_0 = \mathcal{P} \mathcal{L} \mathbf{z}(t)$. Denote $\phi(t, \mathbf{z}) := e^{t\mathcal{L}} \mathcal{Q} \mathcal{L} \mathbf{z}$, then we simplify (36) into

$$\frac{\partial}{\partial t} e^{t\mathcal{L}} \mathbf{z}_0 = \mathcal{P} \mathcal{L} \mathbf{z}(t) + \phi(t, \mathbf{z}_0) + \int_0^t \mathbf{k}(t-s, \mathbf{z}(s)) ds, \quad (37)$$

where $\mathbf{k}(t, \mathbf{z}) := \mathcal{P}\mathcal{L}\phi(t, \mathbf{z}) = \mathcal{P}\mathcal{L}e^{t\mathcal{L}}\mathcal{Q}\mathcal{L}\mathbf{z}$.

Now consider the evolution of $\phi(t, \mathbf{z})$, which is given by

$$\partial_t \phi(t, \mathbf{z}_0) = \mathcal{Q}\mathcal{L}\phi(t, \mathbf{z}_0), \quad (38)$$

with initial condition $\phi(0, \mathbf{z}_0) = \mathcal{Q}\mathcal{L}\mathbf{z}_0 = \mathcal{L}\mathbf{z}_0 - \mathcal{P}\mathcal{L}\mathbf{z}_0 = \bar{\mathbf{f}}(\mathbf{x}_0, \mathbf{e}_0) - \bar{\mathbf{f}}(\mathbf{x}_0, \mathbf{0}) = \mathbf{0}$ since $\mathbf{e}_0 = \mathbf{0}$. Applying \mathcal{P} on both sides of (38) yields

$$\partial_t \mathcal{P}\phi(t, \mathbf{z}_0) = \mathcal{P}\mathcal{Q}\mathcal{L}\phi(t, \mathbf{z}_0) = \mathbf{0},$$

with initial $\mathcal{P}\phi(0, \mathbf{z}_0) = \mathbf{0}$. This implies that $\mathcal{P}\phi(t, \mathbf{z}_0) = \mathbf{0}$ for all t . Hence, applying \mathcal{P} to both sides of (36) yields

$$\frac{\partial}{\partial t} \mathcal{P}\mathbf{z}(t) = \frac{\partial}{\partial t} \mathcal{P}e^{t\mathcal{L}}\mathbf{z}_0 = \mathcal{P}\mathcal{L}\mathbf{z}(t) + \int_0^t \mathcal{P}\mathbf{k}(t-s, \mathbf{z}(s)) ds. \quad (39)$$

Restricting to the first n components, $\mathcal{P}\mathbf{z}(t)$ reduces to $\mathbf{x}(t)$ and $\mathcal{P}\mathbf{k}(t-s, \mathbf{z}(s))$ reduces to $\mathbf{k}(t-s, \mathbf{x}(s))$. Recalling that $\mathcal{P}\mathcal{L}\mathbf{z}(t) = \mathcal{P}\bar{\mathbf{f}}(\mathbf{z}(t)) = \bar{\mathbf{f}}(\mathbf{x}(t), \mathbf{0}) = \mathbf{f}(\mathbf{x}(t))$ completes the proof. \square

B.3 Proof of Theorem 3

Proof. From the definition of $\mathbf{h}(t)$ in (40), we obtain

$$\mathbf{h} = \int_0^t \mathbf{K}(t-s; \mathbf{w})\mathbf{x}(s) ds = \int_{-\infty}^t \mathbf{K}(t-s; \mathbf{w})\mathbf{x}(s) ds = \int_0^\infty \mathbf{K}(s; \mathbf{w})\mathbf{x}(t-s) ds \quad (40)$$

where we used the fact that $\mathbf{x}(t) = \mathbf{0}$ for $t < 0$. Taking derivative on both sides of (40) yields

$$\begin{aligned} \mathbf{h}' &= \int_0^\infty \mathbf{K}(s; \mathbf{w})\mathbf{x}'(t-s) ds = \int_0^\infty \mathbf{K}(s; \mathbf{w})\tilde{\mathbf{f}}(\mathbf{x}(t-s), \mathbf{h}(t-s); \mathbf{A}, \boldsymbol{\eta}) ds \\ &= \int_{-\infty}^t \mathbf{K}(t-s; \mathbf{w})\tilde{\mathbf{f}}(\mathbf{x}(s), \mathbf{h}(s); \mathbf{A}, \boldsymbol{\eta}) ds = \int_0^t \mathbf{K}(t-s; \mathbf{w})\tilde{\mathbf{f}}(\mathbf{x}(s), \mathbf{h}(s); \mathbf{A}, \boldsymbol{\eta}) ds \end{aligned}$$

where we used the fact that $\mathbf{x}'(t) = \tilde{\mathbf{f}}(\mathbf{x}(t), \mathbf{h}(t); \mathbf{A}, \boldsymbol{\eta}) = \mathbf{0}$ for $t < 0$ in the last equality.

If $\mathbf{K}(t; \mathbf{w}) = \sum_l \mathbf{B}_l e^{-\mathbf{C}_l t}$, then we can take derivative of (40) and obtain

$$\begin{aligned} \mathbf{h}'(t) &= \sum_{l=1}^L \frac{d}{dt} \left(\int_{-\infty}^t \mathbf{B}_l e^{-\mathbf{C}_l t} \mathbf{x}(s) ds \right) = \sum_{l=1}^L \left(\mathbf{B}_l \mathbf{x}(t) - \int_{-\infty}^t \mathbf{B}_l \mathbf{C}_l e^{-\mathbf{C}_l t} \mathbf{x}(s) ds \right) \\ &= \sum_{l=1}^L \left(\mathbf{B}_l \mathbf{x}(t) - \mathbf{C}_l \int_{-\infty}^t \mathbf{B}_l e^{-\mathbf{C}_l t} \mathbf{x}(s) ds \right) = \sum_{l=1}^L (\mathbf{B}_l \mathbf{x}(t) - \mathbf{C}_l \mathbf{h}(t)). \end{aligned}$$

Time discretization (14) can then be obtained by finite difference in time with normalized step size 1 and proper scaling of the network parameters $\boldsymbol{\theta}$. \square

B.4 Proof of Theorem 4

Proof. We consider the augmented state $\boldsymbol{\xi}$ and nonlinear dynamics $\bar{\mathbf{g}}(\cdot; \boldsymbol{\theta})$ associated with \mathbf{m} and $\mathbf{g}(\cdot; \boldsymbol{\theta})$, defined as follows:

$$\boldsymbol{\xi}_0 = \begin{bmatrix} \mathbf{m}_0 \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{bmatrix}, \quad \boldsymbol{\xi}_1 = \bar{\mathbf{g}}(\boldsymbol{\xi}_0; \boldsymbol{\theta}) := \begin{bmatrix} \mathbf{g}^1(\mathbf{m}_0; \boldsymbol{\theta}) \\ \mathbf{g}^2(\mathbf{m}_0; \boldsymbol{\theta}) \\ \vdots \\ \mathbf{g}^T(\mathbf{m}_0; \boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \vdots \\ \mathbf{m}_T \end{bmatrix}, \quad (41)$$

where \mathbf{g}^t stands for the composition of $\mathbf{g}(\cdot; \boldsymbol{\theta})$ for t times.

Without overloading the notations, we reuse \mathcal{J} and ℓ of the objective function (18a) and loss function (17) of \mathbf{m} respectively for the augmented state $\boldsymbol{\xi}$. In addition, following [32], we further simplify the notation by combining the K training data into a single variable $\hat{\mathbf{x}} := [\hat{\mathbf{x}}^{(1)}, \dots, \hat{\mathbf{x}}^{(K)}]$; similar

for the state variable \boldsymbol{x} . In this case, the dynamics \boldsymbol{g} is applied to each column of \boldsymbol{x} , and the loss function ℓ is to be interpreted as the average loss as in (17). Furthermore, we temporarily assume the regularization $r(\boldsymbol{\theta}) = 0$ as it is simple to append $\boldsymbol{\theta}$ to the state $\boldsymbol{\xi}$ and merge $r(\boldsymbol{\theta})$ into the loss function $\ell(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}})$. Then the optimal control problem (18) is rewritten as

$$\min_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) := \ell(\boldsymbol{\xi}, \hat{\boldsymbol{\xi}}) + r(\boldsymbol{\theta}) \quad (42a)$$

$$\text{s.t. } \boldsymbol{\xi}_1 = \bar{\boldsymbol{g}}(\boldsymbol{\xi}_0; \boldsymbol{\theta}), \quad \boldsymbol{\xi}_0 = [\boldsymbol{m}_0; \mathbf{0}; \dots; \mathbf{0}]. \quad (42b)$$

Note that (42) is a one-step optimal control with $\bar{\boldsymbol{g}}(\cdot; \boldsymbol{\theta})$. Now by the discrete Pontryagin's Maximum Principle [2], for the state $\boldsymbol{\xi}^*$ optimally controlled by $\boldsymbol{\theta}^*$, there exists a co-state $\boldsymbol{\psi}^*$, such that $\boldsymbol{\xi}^*$ and $\boldsymbol{\psi}^*$ satisfy the following forward and backward equations for $\boldsymbol{\theta} = \boldsymbol{\theta}^*$:

$$\boldsymbol{\xi}_1^* = \bar{\boldsymbol{g}}(\boldsymbol{\xi}_0^*; \boldsymbol{\theta}^*), \quad \boldsymbol{\xi}_0^* = [\boldsymbol{m}_0; \mathbf{0}; \dots; \mathbf{0}], \quad (43a)$$

$$\boldsymbol{\psi}_0^* = \boldsymbol{\psi}_1^* \cdot \nabla_{\boldsymbol{\xi}} \bar{\boldsymbol{g}}(\boldsymbol{\xi}_1^*; \boldsymbol{\theta}^*), \quad \boldsymbol{\psi}_1^* = -\nabla_{\boldsymbol{\xi}} \ell(\boldsymbol{\xi}_1^*, \hat{\boldsymbol{\xi}}), \quad (43b)$$

where

$$\boldsymbol{\xi}_1^* = [\boldsymbol{m}_1^*; \dots; \boldsymbol{m}_T^*] \quad \text{and} \quad \boldsymbol{\psi}_1^* = [\partial_{\boldsymbol{m}_1} \ell(\boldsymbol{\xi}_1^*, \hat{\boldsymbol{\xi}}); \dots; \partial_{\boldsymbol{m}_T} \ell(\boldsymbol{\xi}_1^*, \hat{\boldsymbol{\xi}})] = [\boldsymbol{p}_1^*; \dots; \boldsymbol{p}_T^*]. \quad (44)$$

In addition, $\boldsymbol{\theta}^*$ maximizes the Hamiltonian \mathcal{H} associated with (43):

$$\mathcal{H}(\boldsymbol{\xi}^*, \boldsymbol{\psi}^*; \boldsymbol{\theta}^*) \geq \mathcal{H}(\boldsymbol{\xi}^*, \boldsymbol{\psi}^*; \boldsymbol{\theta}), \quad \forall \boldsymbol{\theta}, \quad \text{where} \quad \mathcal{H}(\boldsymbol{\xi}, \boldsymbol{\psi}; \boldsymbol{\theta}) := \boldsymbol{\psi}_1 \cdot \bar{\boldsymbol{g}}(\boldsymbol{\xi}_0; \boldsymbol{\theta}) - r(\boldsymbol{\theta}). \quad (45)$$

Combining (44), (45), and the definition of H in (19) yields the maximization of total Hamiltonian at the optimal control $\boldsymbol{\theta}^*$:

$$\sum_{t=0}^{T-1} H(\boldsymbol{m}_t^*, \boldsymbol{p}_{t+1}^*; \boldsymbol{\theta}^*) \geq \sum_{t=0}^{T-1} H(\boldsymbol{m}_t^*, \boldsymbol{p}_{t+1}^*; \boldsymbol{\theta}), \quad \forall \boldsymbol{\theta}.$$

For any control $\boldsymbol{\theta}$ and its state and co-state variables $\boldsymbol{\xi}^\theta$ and $\boldsymbol{\psi}^\theta$ following (43) with $\boldsymbol{\theta}$ (also corresponding to \boldsymbol{m}_t^θ and \boldsymbol{p}_t^θ for $t = 0, \dots, T$), we have

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\xi}} \ell(\boldsymbol{\xi}_1^\theta, \hat{\boldsymbol{\xi}}) \cdot \nabla_{\boldsymbol{\theta}} \boldsymbol{\xi}_1^\theta + \nabla_{\boldsymbol{\theta}} r(\boldsymbol{\theta}) \\ &= [\partial_{\boldsymbol{m}_1} \ell(\boldsymbol{\xi}_1^\theta, \hat{\boldsymbol{\xi}}); \dots; \partial_{\boldsymbol{m}_T} \ell(\boldsymbol{\xi}_1^\theta, \hat{\boldsymbol{\xi}})] \cdot [\partial_{\boldsymbol{\theta}} \boldsymbol{g}(\boldsymbol{m}_0^\theta; \boldsymbol{\theta}); \dots; \partial_{\boldsymbol{\theta}} \boldsymbol{g}(\boldsymbol{m}_{T-1}^\theta; \boldsymbol{\theta})] + \nabla_{\boldsymbol{\theta}} r(\boldsymbol{\theta}) \\ &= -\sum_{t=1}^T \left(\boldsymbol{p}_t^\theta \cdot \partial_{\boldsymbol{\theta}} \boldsymbol{g}(\boldsymbol{m}_t^\theta; \boldsymbol{\theta}) + \frac{1}{T} \nabla_{\boldsymbol{\theta}} r(\boldsymbol{\theta}) \right) \\ &= -\sum_{t=1}^T \partial_{\boldsymbol{\theta}} H(\boldsymbol{m}_t^\theta, \boldsymbol{p}_{t+1}^\theta; \boldsymbol{\theta}), \end{aligned}$$

which completes the proof. \square

C Additional Related Work

Network structure inference Inference of diffusion network structure is an important problem closely related to influence estimation. In particular, if the network structure and infection rates are unknown, one often needs to first infer such information from a training dataset of sampled cascades, each of which tracks a series of infection times and locations on the network. Existing methods have been proposed to infer network connectivity [18, 45, 33, 14] and also the infection rates between nodes [37, 17, 19]. Submodular optimization is applied to infer network connectivity [18, 45, 33] by considering the most probable [18] or all [45, 33] directed trees supported by each cascade. One of the early works that incorporate spatio-temporal factors into network inference is introduced in [33]. Utilizing convex optimization, transmission functions [14], the prior probability [37], and the transmission rate [17] over edges are inferred from cascades. In addition to static networks, the infection rates are considered but also in the unobserved dynamic network changing over time [19]. Besides cascades, other features of dynamical processes on networks have been used to infer the diffusion network structures. To avoid using predefined transmission models, the statistical difference of the infection time intervals between nodes in the same cascade versus those not in any cascade was considered in [46]. A given time series of the epidemic prevalence, i.e., the average fraction of infected nodes was applied to discover the underlying network. The recurrent cascading behavior is also explained by integrating a feature vector describing the additional features [50]. A graph signal processing (GSP) approach is developed to infer graph structure from dynamics on networks [35, 11].

D Experiment Supplements

D.1 Implementation details

In our NMF implementation, we use a standard LSTM architecture and 3 dense layers for the RNN ϵ at each time t . Regularization terms using l_1 -norm of all parameters are added to the loss function to promote their sparsity and robustness. Specifically, we use 0.001 to weight A and 0.0001 to all other trainable parameters, respectively. The NMF networks are trained and tested in TensorFlow [1] using Adam optimizer with default parameters ($\text{lr}=0.001$, $\beta_1=0.9$, $\beta_2=0.999$, $\epsilon=1e-8$) on a Linux workstation with Intel 8-Core Turbo 5GHz CPU, 64GB of memory, and an Nvidia RTX 2080Ti GPU. The LSTM model is trained and tested in the same setting as NMF except a fixed regularization weight 0.001 for all trainable parameters. InfluLearner is trained in Matlab, and the number of features is set to 128. All experiments are performed on the same machine. Given ground truth node infection probability x^* , the Mean Absolute Error (MAE) of influence (Inf) and infection probability (Prob) of estimated x are defined by $|\mathbf{1} \cdot (x_t - x_t^*)|$ and $\|x_t - x_t^*\|_1/n$ for every t , respectively.

D.2 Inference of node interdependencies

Due to its highly interpretable structure, NMF can also learn the node inter-dependencies through A . In addition to the quantitative evaluations provided in Section 4, we show the visual appearance of A inferred by NMF in Figure 3. The ground truth A^* and A inferred by NETRATE are also provided for comparison. As we can see, A inferred by NMF is much more faithful to A^* than that by NETRATE. Note that NETRATE requires knowledge of specific diffusion model type (Rayleigh in this test) whereas NMF does not. This result shows that NMF is versatile and robust when only cascade data are available.

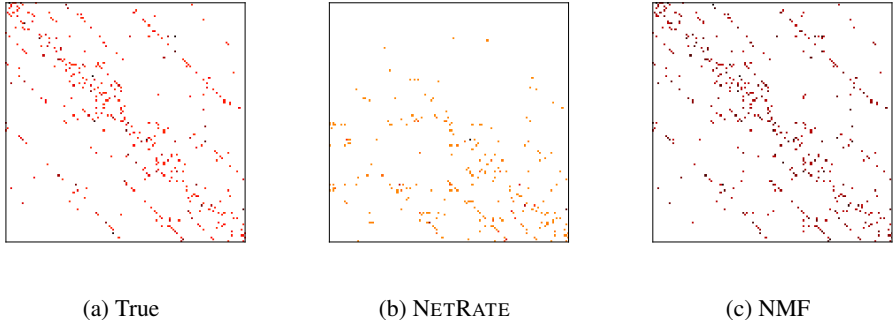


Figure 3: Ground truth A^* (left) and A inferred NETRATE (middle) and NMF (right) under the same color scale using cascaded data from a Hierarchical network with Rayleigh diffusion model.

D.3 Accuracy and Scalability

Accuracy for networks of increasing sizes We test NMF on networks of increasing sizes up to $n=2,048$ with $|\mathcal{E}| = 2n$ for each n using Hierarchical network and exponential diffusion model on cascade data containing 10,000 cascades. We also generate 100 extra cascades with 20%-validation and 80%-test. Figure 4 (a)–(b) shows the MAE of influence (Inf) and infection probability (Prob) estimated by NMF versus time for varying n , which indicate that the error remains low for large networks.

Scalability We compare NMF to InfluLearner in terms of runtime for the influence estimation. For InfluLearner, we draw 200 features. For NMF, the batch size of training cascade data is set to 50 for the network with more than 2,048 nodes, and is 100 for smaller networks. The training is terminated when the average MAE of infection probability on validation data does not decrease for 20 epochs. Figure 4 (c) shows the comparison on runtime (in seconds) of training as we increase the network size n in InfluLearner and NMF. Note that the original implementation of InfluLearner [12] is in

Matlab and the computational time increases drastically in network density, whereas our method retains similar runtime regardless of network density.

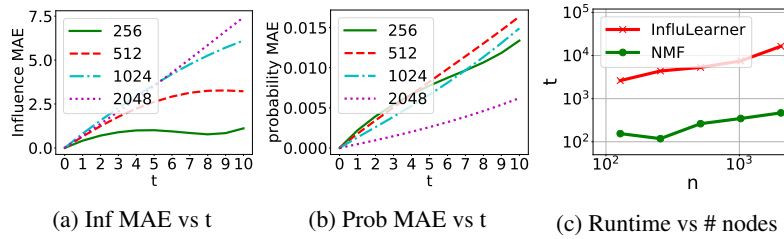


Figure 4: (a)–(b) MAE of influence (Inf) and infection probability (Prob) estimated by NMF for Hierarchical networks with increasing network sizes from 256 to 2048. (c) runtime (in seconds) for training versus network sizes in log-log scale.

D.4 Additional results of infection probability estimation

We test a total of 9 combinations of network structures and diffusion models. Specifically, we generate Hierarchical (Hier), Core-periphery (Core), and Random (Rand) networks, and use Exponential (Exp), Rayleigh (Ray) and Weibull (Wbl) diffusion models on each of these networks. All scale and shape parameters are drawn from $\text{Unif}[0.1, 1]$ and $\text{Unif}[1, 10]$, respectively. Here we stretch NMF and apply to Weibull diffusion model even it has two parameters for each edge. The experiment setting and evaluation metrics are the same as in Section 4. The MAE of influence and node infection probabilities are shown in Figure 5, which shows that NMF consistently performs well with low estimation error after trained by cascade data. Again, it is worth noting that InFLuLearner requires the identity of source node for every infection in the entire cascade during training, which is generally not available in practice nor needed in NMF.

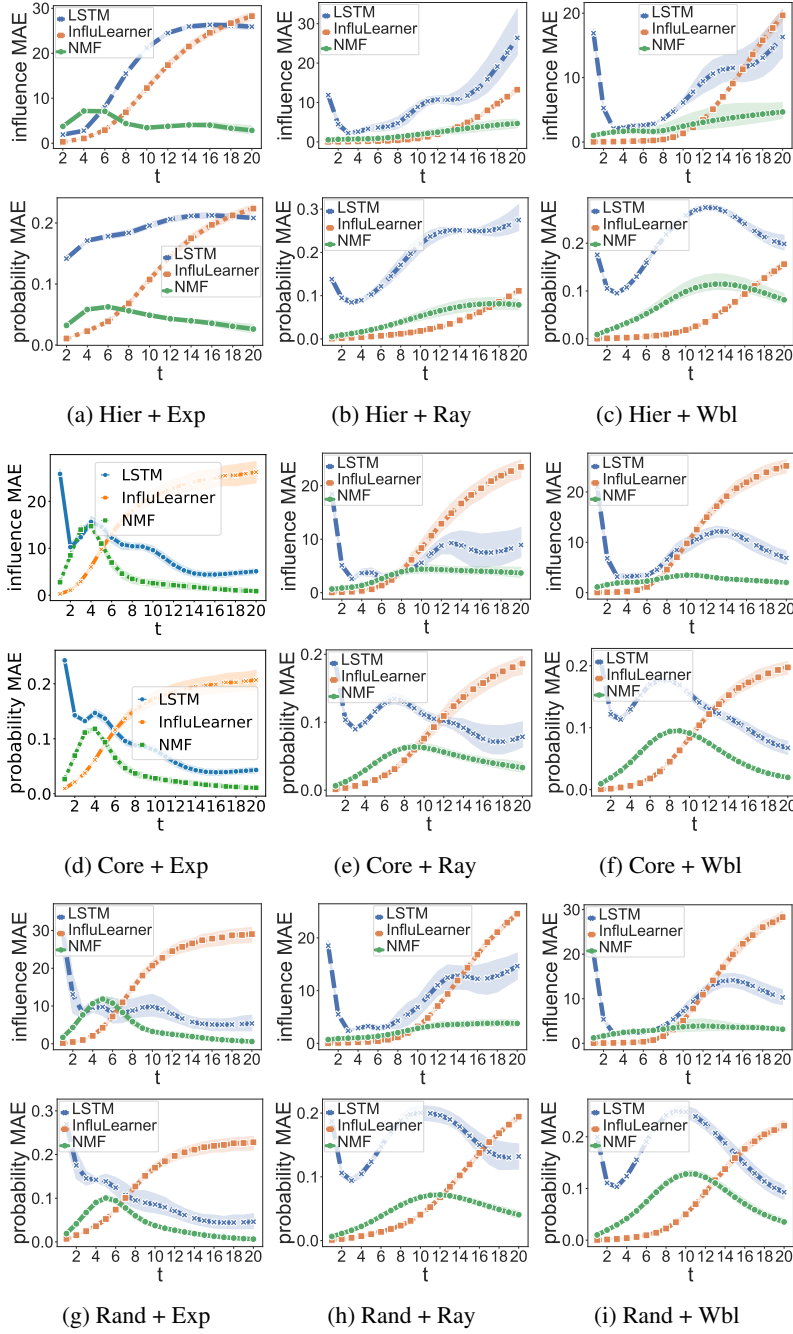


Figure 5: MAE of influence (top) and node infection probability (bottom) by LSTM, InfluLearner, and NMF on each of the 9 different combinations of Hierarchical (Hier), Core-periphery (Core) and Random (Rand) networks, and exponential (Exp), Rayleigh (Ray) and Weibull (Wbl) diffusion models. Mean (centerline) and standard deviation (shade) over 100 tests are shown.