

# 1 Markov State Modeling

Markov State Modeling (MSM) aims to map the slow dynamics of a complex system to an underlying discrete *markovian* process. This involves several steps, including phase space discretization, time discretization and dimensional reduction. There is numerous literature about the mentioned topics [?, ?, ?] and they are not part of this thesis. The analysed models are phenomenological and constructed such that space discretization is performed manually and dimensional reduction is not needed.

## 1.1 Transfer Operator

All existing realisations of a process are combined in an ensemble of trajectories. Considering each trajectory occurs with equal probability, a probability distribution  $P(x)$  can be constructed. This distribution is evolved in time according to a generator, containing all dynamical information by knowledge of all trajectories passing point  $x$ . The generator  $\mathcal{L}$  is denoted by

$$P(x_{t+\tau}) = \mathcal{L} \circ P(x) \quad (1)$$

The initial probability distribution can be chosen according to physical conditions and evolves towards a steady state distribution  $\pi$ . The infinitesimal generator  $\mathcal{L}$  is inconvenient for computational purposes and is replaced by a transfer operator  $\mathcal{T}$  that evolves the system over a time range  $\tau$ . We assume the generator to be *Markovian*, as will be described in detail in section ???. The generator is defined by

$$P_{t+\tau}(y) = \mathcal{T}(\tau) \circ P_t(y) = \int_{\Omega} dx p(x, y; \tau) P_t(x) \quad (2)$$

where  $\Omega$  represents the full phase space of the system and  $p(x, y; \tau)$  is a transition probability in continuous space from state  $y$  to state  $x$  within time  $\tau$ . The Operators are related by

$$\mathcal{T}(\tau) = \exp(\mathcal{L}\tau). \quad (3)$$

where the Eigenvectors  $\Psi_i$  of both are the same and the Eigenvalue  $\lambda_i$  of  $\mathcal{T}$  and  $\Lambda_i$  of  $\mathcal{L}$  are related by  $\lambda = \exp(\tau\Lambda)$ . This relations are needed to construct the infinitesimal generator from a transfer operator. The Eigenvalue decomposition has a maximal Eigenvalue 1 with the corresponding Eigenvector being the stationary distribution. It is assumed to be singular such that the stationary distribution of a system is unique. The smaller Eigenvalues are related to the dynamics of the system as described in section ???.

## 1.2 Markov Property

Time discretization needs special attention such that the resulting process fulfills the *Markov property*.

A Markov process is defined as a *stochastic process* (see section ??) where for n successive timepoints ( $t_1 < t_2 < \dots < t_n$ ) the probability distribution  $P(y_i, t_i)$  given the value  $y_i$  only depends on the information at the previous timepoint  $t_{i-1}$ . We denote

$$P(t_n, y_n | t_1, y_1; t_2, y_2; \dots; t_{n-1}, y_{n-1}) = P(t_n, y_n | t_{n-1}, y_{n-1}) \quad (4)$$

as the time dependent markovian *transition probability* and note that the time evolution of a probability distribution can be expressed by the hierarchy

$$P(t_1, y_1; t_2, y_2; t_3, y_3) = P(t_1, y_1) P(t_2, y_2 | t_1, y_1) P(t_3, y_3 | t_2, y_2), \quad (5)$$

where  $P(t_1, y_1)$  is the initial probability distribution. The Markov process is fully characterized by knowledge of the transition probabilities and the initial probability distribution. Integrating over  $y_2$  and dividing by  $P(t_1, y_1)$  gives

$$P(t_1, y_1 | t_3, y_3) = \int dy_2 P(t_2, y_2 | t_1, y_1) P(t_3, y_3 | t_2, y_2) \quad (6)$$

where the definition of conditional probabilities was used on the left-hand side [?]. Equation ?? is known as the Chapman-Kolmogorov equation and will be used to test if the Markov property is fulfilled.

### 1.2.1 Validate Markovian Process

The time-dependence of the transition probability is dropped to suffice the focus on NESS and a constant timestep length called *lagtime*  $\tau$  is assumed. The time-independent Chapman-Kolmogorov equation becomes

$$P(y_i | y_j, n\tau) = \int dy_1 \int dy_2 \dots \int dy_{j-1} P(y_1 | y_i, \tau) P(y_2 | y_1, \tau) \dots P(y_j | y_{j-1}, \tau) \quad (7)$$

for n time steps and the transition probabilities depend on the lagtime. The state space of a MSM is discretised so we denote  $P(y_i | y_j, \tau) = p_{ij}(\tau)$  as the jump probability from state i to j within time  $\tau$ .

In practice it is cumbersome to check the Chapman-Kolmogorov equation for each element so Prinz et al. [?] suggested to expand the irreducible transition probability matrix  $p_{ij}(\tau)$  in its Eigenvalue decomposition

$$p\Phi_k = \lambda_k \Phi_k. \quad (8)$$

The time evolution of the probability distribution can then be described by

$$P(t) = \sum_k c_k \Phi_k \exp(-\lambda_k t), \quad (9)$$

where  $c_k$  are constants that are determined by the initial state of the system. If the dynamics of the system fulfil detailed balance (i.e. system is in equilibrium), the Eigenvalues and Eigenvectors are real and one can construct a hierarchy

starting with the slowest process  $\lambda_0 = 1 > \lambda_1 > \dots > 0$ . Typically the slow dynamical processes are of interest and the fast processes with low  $\lambda$  are cut off from the expansion. Faster vibrations were already deleted by state space discretisation. The characteristic timescales are identified from equation ?? by  $t_i = \frac{-1}{\log \lambda_i}$  and are calculated for different lagtimes. Figure ?? shows an example of the slowest timescales for varying lagtimes. The region where  $t_i < \tau$  is forbidden because a observed timescale cannot be smaller than the minimal timestep of the Markov process. The timescales reach a plateau for large enough lagtimes, indicating that the Chapman-Kolmogorov equation ?? is valid in this region, i.e. the slowest dynamics of the system are equally described by all choices of lagtimes. A test of consistency of the Eigenvectors is needed for full confirmation of Markovianity, but is as inefficient as comparing single transition probability matrix entries. It is merely a tool to choose a lagtime for further analysis. It is suggested to identify metastable states and compare detailed relaxation probabilities from this state to the trajectory data. The relaxation profile of the MSM should be within the errorbars of the trajectory data as shown in figure ?? for full confirmation, that the model defined by space discretisation and lagtime is markovian for the slowest processes. The Eigenvalue decomposition allows detailed analysis of a MSM by isolating each process and showing detailed probability fluxes involved via the corresponding Eigenvector. The described methods rely on the transition probability matrix being symmetric and the Eigenvalue decomposition being real-valued. In NESS this condition is not met and the Eigenvalues may become complex. A timescale separation with a hierarchy to delete fast processes with small real part of the Eigenvalues is not possible [?]. A similarly powerful tool for analysis of NESS is not known yet, however the Schur-decomposition might be a good candidate for timescale separation. In this thesis, a reference equilibrium system is used to choose a lagtime for further analysis. First-passage-time-distributions (FPTD) are used to confirm markovianity for systems in and out of equilibrium and for analysis of the processes involved.

### 1.3 First-Passage-Time Distribution

First-Passage-Times Distributions are widely used to characterize processes in biology, chemistry and physics and are often associated with a free-energy barrier a system has to overcome. The FPTD contains detailed transition information by collecting numerous realisations of a process. In experiment and simulation of rare events, only the mean of the process is given due to limited observed process realisations [?, ?]. Given a MSM, the metastable states are identified by the steady state distribution (see chapter ??) and the FPTD between all metastable states are calculated. A metastable state can consist of several microstates, the collection of initial states is denoted by  $I$ , of final states by  $F$ . For the purpose of calculating the FPTD from  $I$  to  $F$  the MSM is modified such that all final microstates  $fF$  become a sink, i.e. all jumps out of the metastable state have probability 0 and staying in the state has probability

1

$$\begin{aligned}\tilde{p}_{fj} &= 0 \forall j, \forall f \in F \\ \tilde{p}_{ff} &= 1 \forall f \in F.\end{aligned}\tag{10}$$

An initial state is defined, where full probability is in one microstate  $i \in I$  of the starting metastable state

$$\begin{aligned}\rho_i^{(0)} &= 1 \\ \rho_j^{(0)} &= 0 \forall j \neq i,\end{aligned}\tag{11}$$

where the superscript denotes the number of iterations. The probability distribution is iterated with the modified Markov model  $\hat{p}$  until all probability is trapped in the sink.

$$\vec{\rho}^{(t+1)} = \hat{p} \vec{\rho}^{(t)}.\tag{12}$$

The first-passage probability  $\text{FPT}(t)$  at step  $t$  is then the probability flow in the sink

$$\text{FPT}_{i \rightarrow F}(t) = \sum_{l \in F} \rho_l^{(t)} - \rho_l^{(t-1)}.\tag{13}$$

This calculation is repeated for all initial microstates  $i$ . The FPTD is then calculated by weighting each distribution with the probability of a trajectory starting in each initial microstate

$$\text{FPTD}_{I \rightarrow F}(t) = \sum_{i \in I} \frac{\sum_{k \notin I} \pi_k p_{ki}}{\sum_{i \in I} (\sum_{k \notin I} \pi_k p_{ki})} \text{FPT}_{i \rightarrow F}(t).\tag{14}$$

Knowing the FPTD, all moments of the distribution can be calculated

$$M_{I \rightarrow F}^n = \sum_t \text{FPTD}_{I \rightarrow F}(t) * t^n\tag{15}$$

## 1.4 Entropy Production