



Masterarbeit am Institut für Mathematik der Freien Universität Berlin

Computing the minimal rebinding effect for nonreversible processes

Susanne Röhl

Matrikelnummer: 4364172

susanne.roehl@fu-berlin.de

Betreuer: PD Dr. Marcus Weber

Berlin, 09.04.2017

Contents

Introduction	1
1 Markov State Models	3
1.1 Markov Process	3
1.2 Transfer Operator	7
1.3 Galerkin Projection	11
1.4 Recrossing Effect	17
2 Dominant Structures	22
2.1 Metastability	22
2.2 Spectral Approach	25
2.3 Fuzzy Clustering	28
2.4 Dominant Cycles	32
3 Rebinding Effect in a Given Kinetics	36
3.1 Receptor-Ligand System	36
3.2 Molecular Kinetics as a Projection	42
3.3 Minimizing the Rebinding Effect	47
3.4 Approach for nonreversible processes	52
4 Illustrative Examples	54
4.1 Transition Network Graph	54
4.2 Artificial (bivalent) binding Process	54
Conclusion	55
Bibliography	ii

Introduction

Thesis Structure

Chapter 1 - Markov State Models

We give a short overview of relevant tools from measure theory and Markov processes on measurable state spaces, in order to introduce two different, but strongly related, transfer operators. We show that they propagate probability densities respectively sets(?) of a given Markov process. Then we show how this continuous operator can be projected onto a finite-dimensional space with the aid of a Galerkin discretization. Finally, we analyze the discretization error the possible loss of the Markov property which can occur by this projection.

Chapter 2 - Dominant Structures

We introduce the important concept of metastability as dominant structures of Markov processes. We define metastable sets mathematically and show their relevance for molecular dynamic systems. Furthermore, we reveal their relation to the spectrum of the transfer operator and see how the “best” metastable decomposition can be achieved. At the end of this chapter, we study metastable cycles. They are a different dominant structure which can occur in nonreversible processes.

Chapter 3 - Rebinding Effect

In this chapter, we introduce the crucial point of this thesis, the Rebinding Effect. We describe this effect in the context of a receptor-ligand-system, a special case of a molecular dynamic system, and set this effect in relation to the Recrossing Effect from chapter 1. We apply the methods from chapter 1 and 2 in order rigorously describe a MD system respectively its projection onto a finite subspace. Finally, we compute a minimal bound for the rebinding effect as the solution of an optimization problem.

1 Markov State Models

Stochastic processes, especially Markov processes, are used in many applications in different areas, like biotechnology or Simulations of biomolecular systems (in atomic representation) often require timescales that are far beyond the capacity of computer power currently available (for detailed example see Anton). To get a simulation result in a reasonable time, it makes sense to consider a reduced model of that stochastic process which maintains the relevant dynamical properties while at the same time being less complex. Such reduced models are called “Markov State Models”. There has been a lot of investigations/research activity during the last years, see ...

In order to define/create a Markov State Model, we need at first some basic definitions of stochastic processes, especially Markov processes and how their evolution can be described using the transfer operator. The actual dimension reduction of the process happens by applying a Galerkin projection onto the transfer operator. By that action, states of the original process are clustered/grouped conveniently, such that .. properties/transition rates?.. are preserved..

1.1 Markov Process

We introduce Markov processes, which are a special type of stochastic processes and a generalization of the well-known Markov chains. Markov chains were defined as a memoryless process acting on a finite state space and evolving in discrete time, a behaviour that can be described by a stochastic matrix. For general Markov processes, both these properties can be continuous and thus, we need some more complex formulations and tools in order to describe such processes respectively their time-evolution rigorously.

Transition function

We will denote by $E := (E, \Sigma)$ a *measurable space*, that is a set \mathbb{X} with some σ -algebra Σ defined on it. The triple $\Omega := (\Omega, \mathcal{A}, \mathbb{P})$ will be a *probability space*, that is a measurable space with a probability measure \mathbb{P} defined on it; for detailed information about these basic measure theoretic notations, see Bogachev[1, chapter 1].

A *random variable* $X : \Omega \rightarrow E$ is a *measurable function* from a probability space Ω into a measurable space E , meaning that preimages of measurable sets in E are

measurable in Ω :

$$A \in \Sigma \Rightarrow X^{-1}(A) \in \mathcal{A}.$$

Then the probability measure \mathbb{P} of Ω induces a canonical probability measure on E , by

$$\mu(A) := \mathbb{P}(X \in A) := \mathbb{P}(X^{-1}(A))$$

for all $A \in \Sigma$, see (...).

Definition 1.1. (Stochastic Process)

A family $(X_t)_{t \in \mathbb{T}}$ of random variables $X_t : \Omega \rightarrow E$ on some index set \mathbb{T} is called a *stochastic process* on a state space E .

In the following, we consider stochastic process on real state spaces $E \subset \mathbb{R}^d, d \in \mathbb{N}$, equipped with the Borel- σ -algebra $\Sigma = \mathcal{B}(E)$. In order to introduce Markov processes as a special type of stochastic processes, we need a tool to describe the time evolution or propagation of a process. This can be done using the transition function which describes the propagation of the distribution functions of a stochastic process. ?

Definition 1.2. (Transition function)

A function $p : \mathbb{T} \times E \times \Sigma \rightarrow [0, 1]$ is a *transition function* if it fulfills the following properties:

- i) $x \mapsto p(t, x, A)$ is measurable on E for all $t \in \mathbb{T}$ and $A \in \Sigma$,
- ii) $A \mapsto p(t, x, A)$ is a probability measure for all $t \in \mathbb{T}$ and $x \in E$,
- iii) $p(0, x, E \setminus x) = 0$ for all $x \in E$,
- iv) the Chapman-Kolmogorov equation ?

$$p(t + s, x, A) = \int_E p(t, x, dz) p(s, z, A).$$

holds for all $t, s \in \mathbb{T}, x \in E$ and $A \in \Sigma$.

In this definition, the first three properties just ensure that we get reasonable (measurable) results and that that the process can only be in one state at the same time and not make a jump (a transition in 0-time).

So the transition function $p(t, x, A)$ can be considered as the probability to get into a certain subset A in a time interval t starting from a point x . That follows from the Chapman-Kolmogorov equation, see (...). That means that we can describe the time evolution of a stochastic process by a transition function. In particular, the transition matrix of a Markov chain (time discrete, finite state space) is a special case of the transition function since it fulfills the above properties. why? how?

Markov Process

With the aid of a transition function, we can define Markov processes.

Definition 1.3. (Markov Process)

A stochastic process $(X_t)_{t \in \mathbb{T}}$ on a state space E is a *Markov process* if its transition function fulfills the equation

$$p(t, x, A) = \mathbb{P}(X_{t+s} \in A \mid X_s = x). \quad (1.1)$$

for all $s, t \in \mathbb{T}$, $x \in E$ and $A \in \Sigma$. If that probability is independent from s , then the Markov process is called *time-homogeneous*.

We are especially interested in time-homogeneous processes, which will be presumed from now on. As we can see from the definition, all possible transition probabilities are given and hence, the time evolution of a Markov process is completely described by its transition function. Thus, a Markov process is uniquely determined by its transition function and an initial distribution μ . It is a process that has “no memory” in the sense that only the last known state of the process has an influence on the future of the process, as we can see on the right side of (1.1).

Indeed, there is a one-to-one relation between transition functions and Markov processes, i.e. every homogeneous Markov process defines a transition function and vice versa, see Meyn and Tweedie[20, section 3.4]. The beginning of a Markov process X_t with the transition function p fulfills

$$\mathbb{P}_\mu(X_0 \in A, X_t \in B) = \int_A p(t, x, B) \mu(dx) \quad (1.2)$$

for any $A, B \in \Sigma$, where \mathbb{P}_μ indicates that $X_0 \sim \mu$, or equivalently $\mu(A) = \mathbb{P}(X_0 \in A)$.

The transition function for a Markov process plays the same role as the transition matrix for a Markov chain; it propagates its distributions. If for the transition function we choose $t = 1$ and transitions into one-elementic subsets, then the transition function corresponds to the 1-step transition matrix $[p_{ij}] = T \in \mathbb{R}^{n \times n}$ of a Markov chain. Having introduced the notion of Markov processes, we can now define some important properties and give some examples. ?

Invariant Measure

Definition 1.4. (Invariant measure)

Let $(X_t)_{t \in \mathbb{T}}$ be a Markov process. The probability measure μ is *invariant* w.r.t. $(X_t)_{t \in \mathbb{T}}$ if for all $t \in \mathbb{T}$ and $A \in \Sigma$ we have

$$\int_E p(t, x, A) \mu(dx) = \mu(A).$$

In other words, a measure is invariant wrt a Markov process if the probability to **be** in any subset of the state space is the same as the probability to **get** into that subset by the evolution of the Markov process for any fixed transition time.

also stat.
meas.

Ergodicity

The long-time behaviour of stochastic processes can be described using ergodicity.

Definition 1.5. (ergodic process)

Let $(X_t)_{t \in \mathbb{T}}$ be a Markov process with invariant probability measure μ . Then $(X_t)_{t \in \mathbb{T}}$ is *ergodic* w.r.t. μ if for all functions $u : E \rightarrow \mathbb{R}$ with $\int_E |u| \mu(dx) < \infty$ we have

$$\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T u(X_t) dt = \int_E u(x) \mu(dx).$$

for almost all initial values $X_0 = x_0$.

So a Markov process is ergodic if its time average (left side) is the same as its average over the probability space (right side), known in the field of thermodynamics as its ensemble average. In an ergodic process, the state of the process after a long time is nearly independent of its initial state.

Reversibility

A very useful property of Markov processes is reversibility. A process is reversible if it fulfills the detailed balance equation...; it means that they keep the same probability law even if their movement is considered backwards in time.

Definition 1.6. (reversible process)

Let $(X_t)_{t \in \mathbb{T}}$ be a Markov process with invariant probability measure μ . Then $(X_t)_{t \in \mathbb{T}}$ is *reversible* w.r.t. μ if

$$\int_A p(t, x, B) \mu(dx) = \int_B p(t, x, A) \mu(dx)$$

for all $t \in \mathbb{T}$ and $A, B \in \Sigma$. If μ is unique, then X_t is simply called *reversible*.

If the stochastic transition function is absolutely continuous w.r.t. μ , i.e. ... then reversibility corresponds to $p(t, x, y) = p(t, y, x)$ for all $t \in \mathbb{T}$ and μ -a.e. $x, y \in \mathbb{X}$.

rev. for
Markov
chain bzw.
abs. cont.
meas.
rev. follows
ex. stat./inv.
dist.?
det.bal.?

Example: Markov Chain

Let $(X_t)_{t \in \mathbb{T}}$ be a *Markov chain*, that is a Markov process on discrete time $\mathbb{T} = \mathbb{N}$ and finite state space $E = \{1, \dots, n\}$. Its invariant measure corresponds to the stationary

distribution π . Since we are considering 1-step transitions, the associated transition function is given by $p(x, y) := p(1, x, y)$ and corresponds to the entries of the *transition matrix* $T \in \mathbb{R}^{n \times n}$, that is

$$T_{xy} = p(x, y) = \mathbb{P}(X_1 = y \mid X_0 = x).$$

The propagation of a probability distribution $v_0 \in \mathbb{R}^n$ in the state space can be written as $v_1^T = v_0^T T$, where v_0^T denotes the transposed vector of v_0 . The invariant measure $\pi \in \mathbb{R}^n$ satisfies $\pi^T = \pi^T T$. It is a left eigenvector to the eigenvalue 1 of the transition matrix. If the matrix is irreducible, such an eigenvector exists, due to Perron-Frobenius theorem.

1.2 Transfer Operator

With the previously defined transition function, we have a tool to describe the propagation of **distributions** of stochastic processes. Now we are going to introduce an operator that propagates **probability densities** of Markov processes. Before defining such an operator, we have to specify the space of functions the operator is acting on.

similarly to
trans. matr.
of MC?
resp. sets?

L^r -Spaces

It seems natural to define such a density propagating operator as acting on $L^1(\mu)$, the Banach space that includes all probability densities with respect to μ . However, it is sometimes advantageous to restrict the analysis to $L^2(\mu)$, since this can lead to a self-adjoint operator. As there are different motivations for the choice of such a space, we define an operator which acts on $L^r(\mu)$ -spaces, i.e. spaces of r -integrable functions.

only re-
versible

Definition 1.7. (L^r -Spaces)

Let (E, Σ, μ) a measure space. Then we define the corresponding L^r -spaces as equivalence classes of measurable functions

\mathbb{C} ?

$$L^r(E, \Sigma, \mu) = \{f : E \rightarrow \mathbb{R} \mid \int_E |f(x)|^r \mu(dx) < \infty\}$$

for $1 \leq r < \infty$ and

$$L^\infty(E, \Sigma, \mu) = \{f : E \rightarrow \mathbb{R} \mid \text{ess sup}_{x \in E} |f(x)|^r \mu(dx) < \infty\},$$

with the corresponding norms $\|\cdot\|_r$ and $\|\cdot\|_\infty$.

In these equivalence classes, two functions f, g are identified if $f = g$ μ -almost everywhere, see Werner[41, section I.1]. It is clear from the context, which measure

space (E, Σ, μ) is in consideration, we just write shortly $L^r(\mu) := L^r(E, \Sigma, \mu)$. Due to Hölders inequality, we have $L^r(\mu) \subset L^s(\mu)$ for all $1 \leq s \leq r \leq \infty$. All L^r -spaces are Banach spaces, though $L^2(\mu)$ is the only one which can be equipped with a canonical scalar product and thereby becomes a Hilbert space, see Werner[41, section V.1]. For $f, g \in L^2(\mu)$, the scalar product is defined as

compl.?

$$\langle f, g \rangle_\mu := \int_E f(x) \overline{g(x)} \mu(dx).$$

Now let ν_0 be the density function of a given start distribution. Then the density function of a subset $A \in \Sigma$ at time t is given in terms of the transition function by

$$\nu_t(A) = \int_E \nu_0 p(t, x, A) \mu(dx).$$

On the other hand, the density ν_t is given by

$$\nu_t(A) = \int_A \nu_t(x) \mu(dx).$$

Forward and Backward Transfer Operator

The two above equations result in the following intuitive definition of a transfer operator which should “propagate” probability densities according to a given Markov process. But instead of limiting us to density functions, we define the transfer operator as acting on any r -integrable function.

propagates
subensem-
bles?

Definition 1.8. (Propagator or Forward Transfer Operator)

Let $p : \mathbb{T} \times E \times \Sigma \rightarrow [0, 1]$ be the transition function of a Markov Process $(X_t)_{t \in \mathbb{T}}$ and μ be an invariant measure of $(X_t)_{t \in \mathbb{T}}$. The semigroup of *propagators* or *forward transfer operators* $\mathcal{T}^t : L^r(\mu) \rightarrow L^r(\mu)$ with $t \in \mathbb{T}$ and $1 \leq r \leq \infty$ is defined via

semigroup?

$$\int_A \mathcal{T}^t \nu(y) \mu(dy) = \int_E \nu(x) p(t, x, A) \mu(dx) \quad (1.3)$$

for all $A \in \Sigma$ and $\nu \in L^r(\mu)$.

The propagator is well-defined on the Banach spaces $L^r(\mu)$, $1 \leq r \leq \infty$, see [14]. We will list already some properties of this operator which will be useful in the following chapters. \mathcal{T}^t is a *Markov operator*, i.e. it conserves the norm, $\|\mathcal{T}^t \nu\|_1 = \|\nu\|_1$, and is positive, $\mathcal{T}^t \nu \geq 0$ for $\nu \geq 0$. $\mathcal{T}^t \nu_0$ describes the transport of the function ν_0 in time t by the underlying dynamics given by the process X_t and weighted with respect to μ via

$$\nu_0 \mapsto \nu_t = \mathcal{T}^t \nu_0.$$

Since μ is invariant, we have that the characteristic function of the state space is invariant under the action of \mathcal{T}^t , that is

$$\mathcal{T}^t \mathbb{1}_E = \mathbb{1}_E.$$

?
 $\mathbb{1} := \mathbb{1}_E$?

It means that \mathcal{T}^t has the eigenvalue 1 which corresponds to its eigenfunction $\mathbb{1}_E$.

Definition 1.9. (Backwards Transfer Operator¹)

The *backwards transfer operator* $\mathcal{U}^t : L^r(\mu) \rightarrow L^r(\mu)$ with $t \in \mathbb{T}$ and $1 \leq r \leq \infty$ is defined by

$$\mathcal{U}^t f(x) = \int_E f(y) p(t, x, dy). \quad (1.4)$$

We have again 1 as eigenvalue to the eigenfunction $\mathbb{1}_E$, that is for all $t \in \mathbb{T}$ we have

$$\mathcal{U}^t \mathbb{1}_E = \mathbb{1}_E.$$

$\|\mathcal{U}f\| \leq \|f\|$
 $\|\mathcal{U}\| \leq 1$?

The operator \mathcal{U}^t is *adjoint* to \mathcal{T}^t , denoted by $(\mathcal{T}^t)^* = \mathcal{U}^t$, that is they are related via the duality bracket, namely for all $f \in L^p(\mu), g \in L^q(\mu)$ with $\frac{1}{p} + \frac{1}{q} = 1$, we have

$\mathcal{U} := \mathcal{U}^t$?

$$\langle \mathcal{T}^t f, g \rangle_\mu = \langle f, \mathcal{U}^t g \rangle_\mu.$$

We again remark that both forward as well as backward operator can be defined on arbitrary $L^r(\mu)$ -spaces. But the previous equation shows us that either the choice $p = q = 2$ or the choice $p = 1, q = \infty$ or conversely, make sense, in order that we obtain this useful adjointness/duality-relation of the two operators.

If we compare the equations (1.3) and (1.4), the notion of “forward” and “backwards” becomes clear. For the forward case, the state average with respect to f is taken over all initial states x which are propagated forward in time. In the backward case, we take the state average over all final states y .

If the state space is finite and the corresponding process reversible, then we can see the relation of the forward and backward operator still better. Then the forward operator corresponds to the transition matrix, propagating probability distributions, while the backward operator corresponds to the *transposed* transition matrix, propagating subsets.

?

Spectrum of Transfer operator

Later in this thesis, we will be interested in examining the spectrum of the transfer operator of a given Markov process. The following theorems give us an important insight about the spectrum and its relation to the reversibility of the process.

¹This nomenclature is motivated by the fact that for some models the forward transfer operator is related to the forward Kolmogorov relation, while the backward transfer operator is related to the backward Kolmogorov relation.

Definition 1.10. (Self-adjoint Operator)

An operator \mathcal{T} on $L^2(\mu)$ is called *self-adjoint* if for all $f, g \in L^2(\mu)$ we have

$$\langle f, \mathcal{T}g \rangle_\mu = \langle \mathcal{T}f, g \rangle_\mu.$$

Theorem 1.11. (Werner[41, theorem VI.1.2, theorem VI.1.3, lemma VI.3.1])

Let X be a Banach space and $\mathcal{T} : X \rightarrow X$ a linear continuous operator. Then we have

$$|\lambda| \leq \|\mathcal{T}\| \text{ for all } \lambda \in \sigma(\mathcal{T}).$$

If X is a Hilbert space, then

$$i) \sigma(\mathcal{T}^*) = \{\bar{\lambda} \mid \lambda \in \sigma(\mathcal{T})\}$$

$$ii) \text{ if } \mathcal{T} \text{ is self-adjoint, i.e. if } \mathcal{T}^* = \mathcal{T}, \text{ then } \sigma(\mathcal{T}) \subset \mathbb{R}$$

$$iii) \text{ if } \mathcal{T} \text{ is self-adjoint, then each two eigenfunctions corresponding to different eigenvalues are orthogonal}$$

needed?

Since we know that the operator norm of any transfer operator \mathcal{T} is 1, it follows immediately from theorem 1.11 that its spectrum $\sigma(\mathcal{T})$ is contained in the unit circle of the complex plane, that is we have $|\lambda| \leq 1$ for all $\lambda \in \sigma(\mathcal{T}) \subset \mathbb{C}$.

Theorem 1.12. (Huisinga[14, proposition 1.1])

Let $\mathcal{T}^t : L^2(\mu) \subset L^1(\mu) \rightarrow L^2(\mu)$ be the propagator corresponding to the Markov process $(X_t)_{t \in \mathbb{T}}$. Then \mathcal{T}^t is self-adjoint with respect to the scalar product $\langle \cdot, \cdot \rangle_\mu$ in $L^2(\mu)$ if and only if $(X_t)_{t \in \mathbb{T}}$ is reversible.

Thus, the transfer operator of a reversible process has a spectrum $\sigma(\mathcal{T}) \subset [-1, 1]$. Furthermore, theorem 1.11 guarantees us that the spectrum of a self-adjoint operator is equal to the spectrum of its adjoint. Thus, if we are given a reversible process, it doesn't matter if we examine the spectrum of the forward or the backward transfer operator.

Infinitesimal Generator

For $\mathbb{T} = \mathbb{R}$ the Chapman-Kolmogorov property of the transition functions makes the family $\{\mathcal{T}^t\}_{t \in \mathbb{R}}$ a continuous *semigroup* due to

$$\mathcal{T}^{t+s} = \mathcal{T}^t \mathcal{T}^s.$$

proof
also backw.?

This leads to the following definition of the the infinitesimal generator.

time-indep.?

Definition 1.13. (Infinitesimal Generator)

For the semigroup of propagators or forward transfer operators $\mathcal{T}^t : L^r(\mu) \rightarrow L^r(\mu)$ with $t \in \mathbb{T}$ and $1 \leq r \leq \infty$ we define $\mathcal{D}(L)$ as the set of all $\nu \in L^r(\mu)$ s.t. the strong limit

$$\mathcal{Q}\nu = \lim_{t \rightarrow 0} \frac{\mathcal{T}^t \nu - \nu}{t}$$

exists. Then the operator $\mathcal{Q} : \mathcal{D}(L) \rightarrow L^r(\mu)$ is called the *infinitesimal generator* corresponding to the semigroup \mathcal{T}^t .

The infinitesimal generator is an operator which describes the behaviour of a Markov process in infinitesimal time. That becomes clear by the relation $\mathcal{T}^t = \exp(t\mathcal{Q})$ in $L^2(\mu)$. We can say that \mathcal{Q} “generates” the semigroup of transfer operators since the whole semi-group of transfer operators can be derived from it. ref

Therefore, the eigenvalues $1 = \lambda_1, \dots, \lambda_m$ of the propagator \mathcal{T}^t are related to the eigenvalues $0 = \Lambda_1, \dots, \Lambda_m$ of the generator \mathcal{Q} via not yet discrete?

$$\lambda_k = \exp(t\Lambda_k)$$

for all $1 \leq k \leq m$. Their corresponding (associated) eigenfunctions are identical. Thus, the stationary distribution of \mathcal{T}^t is the solution of $\mathcal{Q}\pi = 0$; $\mathcal{Q}\mathbb{1} = 0$. properties of generator?

1.3 Galerkin Projection

So far we considered Markov processes on very large, possibly continuous, state spaces. For many applications, simulations of a given process are needed in order to obtain informations about the corresponding system. But computations on **large** state spaces require an enormous amount of computation power and time. With larger state spaces, the computation effort increases exponentially fast, see “curse of dimension”. Therefore, we are interested in reducing the number of states in order to make computations more feasible. Such a reducing can for instance be done by “grouping similar states together”. In this section, we will develop a mathematical concept/tool which enables us to create such a reduced model. see ..

Galerkin Projection

The first step in order to create our desired finite process is to determine a convenient finite state space $D \subset L^2(\mu)$. For this purpose, we choose a partition of unity as a basis, which is a generalization of a set of characteristic functions. That more general idea gives us more flexibility for later applications. The relevance of the choice of a partition of unity for the projection will be clarified in section 2.3.

Definition 1.14. (Partition of Unity)

A family of measurable functions $\{\chi_1, \dots, \chi_n\} \subset L^2(\mu)$ is called a *partition of unity* if the following two conditions are fulfilled:

i) The χ_i are non-negative and linear independent

ii) $\sum_{i=1}^n \chi_i(x) = 1$ for all $x \in E$.

$\chi_i : E \rightarrow [0, 1]$?

Definition 1.15. (Galerkin Projection)

Let $\{\chi_1, \dots, \chi_n\}$ be a partition of unity, $D = \text{span}\{\chi_1, \dots, \chi_n\}$ the associated finite-dimensional ansatz space and $\hat{S} \in \mathbb{R}^{n \times n}$ with $\hat{S}_{kj} = \langle \chi_k, \chi_j \rangle_\mu$. The *Galerkin projection* onto D is defined by $G : L^2(\mu) \rightarrow D$ via

$$G\nu = \sum_{k,j=1}^n \hat{S}^{-1}(k, j) \langle \chi_k, \nu \rangle_\mu \chi_j. \quad (1.5)$$

The matrix \hat{S} is invertible since it is the Gramian matrix of linear independent functions. In the easy case that the $\{\chi_1, \dots, \chi_n\}$ are the characteristic functions belonging to a full partition $\{A_1, \dots, A_n\}$, equation (1.5) becomes

“weighted” orthogonal projection?

$$G\nu = \sum_{k=1}^n \frac{1}{\mu(A_k)} \langle \chi_k, \nu \rangle_\mu \chi_k,$$

since the χ_i are orthogonal which means that $\chi_k \chi_j = 1$ if $j = k$ and 0 otherwise.

A Galerkin projection can be applied on the transfer operator of a Markov process.

Definition 1.16. (Projected Transfer Operator)

Let $\mathcal{P} := \mathcal{P}^t$ be the transfer operator of a Markov process on a state space E with unique invariant measure μ , $\{\chi_1, \dots, \chi_n\}$ be a partition of unity and G the Galerkin projection onto the associated subspace D . Then an operator of the form

$$G\mathcal{P}G : L^2(\mu) \rightarrow D$$

is called *projected transfer operator* and we abbreviate it by $G(\mathcal{P})$.

Matrix Representation

Since we are interested in transitions inside of our smaller (projected) space, we want to propagate n -dimensional vectors by the projected transfer operator. For this reason, we consider now the projection of the *restricted* transfer operator $G\mathcal{P}|_D : D \rightarrow D$, which will be denoted $G(\mathcal{P})$ as well.

We remind that every linear map between finite-dimensional vector spaces can be represented by a matrix which is determined by chosen bases. Thus we can write the projected transfer operator as a $n \times n$ -matrix in the following useful way.

GPG same as GP_D ? see .. why linear?

Theorem 1.17. (Sarich [26])

Let \mathcal{P} be the transfer operator of a Markov process, $\{\chi_1, \dots, \chi_n\}$ a partition of unity and $G(\mathcal{P})$ the Galerkin projection of the transfer operator onto the associated subspace. Then $G(\mathcal{P})$ has a matrix representation

$$P_c = TS^{-1},$$

where

$$S_{kj} = \frac{\hat{S}(k, j)}{\langle \chi_k, \mathbb{1} \rangle_\mu} = \frac{\langle \chi_k, \chi_j \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} \quad (1.6)$$

and

$$T_{kj} = \frac{\langle \chi_j, \mathcal{P}\chi_k \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu}. \quad \mathbb{1}_E?$$

Proof. Remember that P_c is a (left) matrix representation of $G(\mathcal{P})$ with respect to a basis $\{\psi_1, \dots, \psi_n\}$ of D if for any function $f : D \rightarrow D$ with

$$f = \sum_{i=1}^n \alpha_i \psi_i \quad \text{and} \quad G(\mathcal{P})f = \sum_{i=1}^n \beta_i \psi_i \quad (1.7)$$

it holds that

$$(\alpha_1, \dots, \alpha_n)P_c = (\beta_1, \dots, \beta_n). \quad (1.8)$$

We assume that (1.7) is true and we aim to show (1.8). For that purpose, we choose a basis $\{\psi_1, \dots, \psi_n\}$ of D with

$$\psi_k = \frac{\chi_k}{\langle \chi_k, \mathbb{1}_E \rangle_\mu}. \quad (1.9)$$

As $G(\mathcal{P})$ is a linear map, we have $G(\mathcal{P})f = \sum \alpha_i G(\mathcal{P})\psi_i$. We exploit this fact, as well as the definitions of the Galerkin projection and the basis to compute

$$\begin{aligned} G(\mathcal{P})f &= \sum_{k=1}^n \alpha_k G(\mathcal{P})\psi_k \\ &\stackrel{(1.5)}{=} \sum_{k,l,j=1}^n \alpha_k \hat{S}^{-1}(j, l) \langle \chi_j, \mathcal{P}\psi_k \rangle_\mu \chi_l \\ &\stackrel{(1.9)}{=} \sum_{k,l,j=1}^n \alpha_k \hat{S}^{-1}(j, l) \langle \chi_j, \mathcal{P}\psi_k \rangle_\mu \langle \chi_l, \mathbb{1}_E \rangle_\mu \psi_l \\ &\stackrel{(1.7)}{=} \sum_{l=1}^n \beta_l \psi_l. \end{aligned}$$

Comparing the coefficients of the last two equations, we can express β_l as

$$\begin{aligned}\beta_l &= \sum_{k,j=1}^n \alpha_k \hat{S}^{-1}(j, l) \langle \chi_l, \mathbb{1}_E \rangle_\mu \langle \chi_j, \mathcal{P}\psi_k \rangle_\mu \\ &= \sum_{k=1}^n \alpha_k \underbrace{\sum_{j=1}^n \hat{S}^{-1}(j, l) \langle \chi_l, \mathbb{1}_E \rangle_\mu \frac{\langle \chi_j, \mathcal{P}\chi_k \rangle_\mu}{\langle \chi_k, \mathbb{1}_E \rangle_\mu}}_{\stackrel{!}{=}(P_c)_{kl}}.\end{aligned}\tag{1.10}$$

The underbraced term **should** be equal to $(P_c)_{kl}$ because we wish that (1.8) is fulfilled. Thus, we compute the (k, l) -th entry of $P_c = TS^{-1}$, employing the fact that $(S^{-1})_{jl} = (\hat{S}^{-1})_{jl} \langle \chi_l, \mathbb{1} \rangle$, as

notation!

$$\begin{aligned}(TS^{-1})_{kl} &= \sum_{j=1}^n T_{kj}(S^{-1})_{jl} \\ &= \sum_{j=1}^n \frac{\langle \chi_j, \mathcal{P}\chi_k \rangle}{\langle \chi_k, \mathbb{1}_E \rangle} (\hat{S}^{-1})_{jl} \langle \chi_l, \mathbb{1}_\mathbb{X} \rangle\end{aligned}$$

and discover that it is equal to the underbraced term in (1.10). Hence, (1.8) is true and therefore P_c is the requested matrix representation of $G(\mathcal{P})$. \square

Theorem 1.18. *The matrices S and T from theorem 1.17 are stochastic.*

non-neg.

Proof. In order to be stochastic, each row must sum up to 1. We exploit the partition of unity property $\sum_j \chi_j = 1$ for all j and the aforementioned properties $\mathcal{P}\mathbb{1} = \mathbb{1}$ and $\mathcal{P}^*\mathbb{1} = \mathbb{1}$ of a transfer operator respectively its adjoint:

$$\begin{aligned}\sum_{j=1}^n S_{kj} &= \frac{\langle \chi_k, \sum_j \chi_j \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} = \frac{\langle \chi_k, \mathbb{1} \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} = 1, \\ \sum_{j=1}^n T_{kj} &= \frac{\langle \sum_j \chi_j, \mathcal{P}\chi_k \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} = \frac{\langle \mathbb{1}, \mathcal{P}\chi_k \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} = \frac{\langle \mathcal{P}^*\mathbb{1}, \chi_k \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} = 1.\end{aligned}$$

\square

Since S and T are both stochastic matrices, they have $\mathbb{1}_D$ as right eigenvector to the eigenvalue 1. It implies that the same holds for P_c , that is its rows sum up to 1 and thus the product TS^{-1} is **at least pseudostochastic**. But nonnegativity is not assured since inverting S can provoke negative entries. The non-negativity depends on the choice of the partition of unity. As we will see, there are examples such that TS^{-1} is a stochastic matrix.

Theorem 1.19. *The matrix representation P_c from theorem 1.17 has the left eigenvector $\hat{\mu} \in D$ with the entries*

$$\hat{\mu}_j = \langle \mathbb{1}, \chi_j \rangle_\mu = \int_E \chi_j(x) \mu(dx).$$

Proof. We observe that $\hat{\mu}^T S = \hat{\mu}^T$ and $\hat{\mu}^T T = \hat{\mu}^T$ since

$$(\hat{\mu}^T S)_j = \sum_{k=1}^n \langle \mathbb{1}, \chi_k \rangle_\mu \frac{\langle \chi_k, \chi_j \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu} = \langle \mathbb{1}, \chi_j \rangle_\mu = \hat{\mu}_j$$

and

$$(\hat{\mu}^T T)_j = \sum_{k=1}^n \langle \chi_j, \mathcal{P} \chi_k \rangle_\mu = \langle \chi_j, \mathcal{P} \mathbb{1} \rangle_\mu = \langle \chi_j, \mathbb{1} \rangle_\mu = \hat{\mu}_j.$$

We can deduce that $\hat{\mu}^T P_c = \hat{\mu}^T T S^{-1} = \hat{\mu}^T S^{-1} = \hat{\mu}^T S S^{-1} = \hat{\mu}^T$. \square

Remark 1.20. We remark that in theorem 1.17 we stated the matrix representation of the *(forward) propagator* of a projected Markov process. As the backward operator is the adjoint operator of the propagator, the matrix representation of its projection is just the transpose $P_c^T = S^{-1}T$. This result and its properties are computed completely analogously to the previous proofs.

nop

P_c propagates probability distributions, while P_c^T propagates sets(?).

Example: Full Partition Discretization

Let A_1, \dots, A_n be a partition of the state space E , i.e. they are pairwise disjoint sets such that $\cup A_i = \mathbb{X}$. We consider the family of the corresponding characteristic functions

$$\chi_i(x) = \mathbb{1}_{A_i}(x).$$

Since they are orthogonal, the matrix S is the identity matrix and therefore, the matrix representation of the Galerkin projection is $P_c = T$. We can compute it by combining

$$\langle \mathcal{P} \mathbb{1}_{A_i}, \mathbb{1}_{A_j} \rangle_\mu = \int_{A_j} (\mathcal{P} \mathbb{1}_{A_i})(x) \mu(dx) \stackrel{(1.3)}{=} \int_{A_i} p(t, x, A_j) \mu(dx) \stackrel{(1.2)}{=} \mathbb{P}_\mu(X_t \in A_j, X_0 \in A_i)$$

and

$$\langle \mathbb{1}_{A_i}, \mathbb{1} \rangle_\mu = \int_{\mathbb{X}} \mathbb{1}_{A_i}(x) \mu(dx) = \mu(A_i) = \mathbb{P}_\mu(X_0 \in A_i).$$

The entries of the resulting matrix representation are given by

$$T_{ij} = \frac{\langle \mathcal{P} \mathbb{1}_{A_i}, \mathbb{1}_{A_j} \rangle_\mu}{\langle \mathbb{1}_{A_i}, \mathbb{1} \rangle_\mu} = \mathbb{P}_\mu(X_t \in A_j \mid X_0 \in A_i).$$

Thus, P_c represents a Markov chain whose state space consists of the partition sets A_i , e.g. each A_i is a “macro state” of the projected process. The stationary distribution of this Markov chain P_c is just the projection of the invariant measure μ onto D .

For a full partition discretization, the matrix S is a diagonal matrix. If we choose a partition of unity that is *close* to a full partition, i.e. we choose *almost characteristic functions*, then the matrix S is not diagonal, but close to that. We will later see the consequences of that fact regarding to the examination of the *rebinding effect*.

Properties of Galerkin Projection

As the matrix representation of a projected transfer operator is in general **not** a stochastic matrix, and stochastic matrices are in a one-to-one relation with Markov chains, we can immediately deduce that the process can lose its Markovianity by projecting it onto a subspace. This possible loss of Markovianity is certainly a really undesirable effect. But before examining that later in section 1.4, let us now first analyze further, hopefully **nice**, properties of the matrix representation P_c .

We already know that the matrices S and T from Theorem 1.17 are stochastic matrices. This leads to some good properties of P_c :

- The eigenvalue $\lambda = 1$ of P_c has the associated right-eigenvector $e = (1, \dots, 1)^T$ and left-eigenvector $\hat{\mu}^T$ from theorem 1.19.
- If \mathcal{P} is self-adjoint in L^2 , then $G(\mathcal{P})$ as well. Then the matrices S and T are self-adjoint with respect to the discrete scalar product

$$\langle u, v \rangle_{\hat{\mu}} = \sum_{i=1}^n u_i v_i \hat{\mu}_i.$$

P_c self-adj.?

$$\begin{aligned} \langle Av, w \rangle &= \\ \langle v, Aw \rangle \end{aligned}$$

Since self-adjointness of the operator is equivalent to reversibility of the corresponding process, see theorem 1.12, detailed balance equation (e.g. $\hat{\mu}_k T_{kl} = \hat{\mu}_l T_{lk}$ for all $k, l = 1, \dots, n$) is fulfilled for both S and T .

thus, eigenv. of S and T real and in $[-1, 1]$

- If the transfer operator has a simple and dominant eigenvalue 1 and the continuous part of the spectrum is bounded away from the discrete part, then the process is irreducible and aperiodic which is inherited by the matrix T . In particular T has the simple and dominant eigenvalue $\lambda = 1$ which is the only eigenvalue with $|\lambda| = 1$ and the discrete invariant density $\hat{\mu}$ is the unique invariant density of T .
- As seen in the last example, a full-partition projection yields the transition matrix $P_c = T$ of a Markov chain describing transitions between the partition sets.

Conclusion, Interpretation

The discretization maintains/inherits many important properties of the original process. Interpretation of S as a “stochastic mass matrix” and T as a “coupling matrix” (transition matrix) ?

Projected infinitesimal generator

The Galerkin projection of an infinitesimal generator yields a similar matrix representation as the transfer operator. It can also be written as the product of two stochastic matrices, one of them being the inverted mass matrix of the partition of unity functions.

Theorem 1.21. *Let $\mathcal{Q} : L^2(\mu) \rightarrow L^2(\mu)$ be a generator of a semigroup of transfer operators with unique invariant measure μ and satisfying $\mathcal{Q}\mathbb{1}_E = 0$. Let χ be a partition of unity with a projection G onto the associated subspace spanned by χ . Then the projected generator $G(\mathcal{Q})$ has the matrix representation $Q = RS^{-1}$ with the stochastic mass matrix S from (1.6) and*

$$R(k, j) = \frac{\langle \chi_j, \mathcal{Q}\chi_k \rangle_\mu}{\langle \chi_k, \mathbb{1} \rangle_\mu}$$

The eigenvalue problem of Q is equivalent to the generalized eigenvalue problem $Ru = \Lambda Su$. For both Q and R the largest eigenvalue is $\lambda = 0$. The associated right eigenvector is $e = (1, \dots, 1)^T$, the associated left eigenvector is $\hat{\mu}^T$ from theorem 1.19.

Proof. The matrix representation of $G\mathcal{Q}G$ can be shown similar to the proof of theorem 1.17. For the other properties see... \square

There are obviously many possibilities/options to make a Galerkin discretization/projection of the propagator/generator of a given process. So far, we gave the example of a full-partition discretization, which results in a very easy matrix representation. But as arbitrary partitions of unity χ_1, \dots, χ_n are allowed for a Galerkin projection, there will be more interesting results. In chapter 2 we are going to see which choice of χ gives us a *good* discretization of our process in the sense that it maintains certain desired properties; in our case the long-time behaviour of the process using so called *metastability*.

1.4 Recrossing Effect

In this section, we are going to examine the so called *Recrossing Effect* which is one of the main topics of this thesis. This effect, occuring when projecting a process onto a

relation to iteration error

smaller state space, may spoil the Markov Property of the process.

We are going to analyze it by means of an easy example. Additionally, we have to face the problem that a projected process/transfer operator does not necessarily behaves/propagates as the original process. We will explain the relevance of this *iteration error* without going into further details, since in the following chapters we will be able to apply a projection where this error/deviation vanishes.

measure by S ?
kind of memory effect

Initial Situation

Assume we are given a Markov process $(X_t)_{t \in \mathbb{T}}$ on a continuous or very large state space E , described by the transfer operator $\mathcal{P} := \mathcal{P}(\tau)$. In order to get a discrete process out of it, we are going to project the time onto \mathbb{N} and the state space onto a finite set $\{1, \dots, n\}$. Discretizing the time can be done naturally without problems since for every lag-time $\tau > 0$, the process $(X_{k\tau})_{k \in \mathbb{N}}$ is again Markovian.

why?

However, the state-space discretization has to be observed a bit more elaborated. Let's do this on the example of a full partition discretization. We consider the operator $G(\mathcal{P}^k) := G\mathcal{P}^k G$, that is we first propagate the process and project it afterwards. Then for all k -multiples of τ , we assign the current state of the original process X_t to the projected process \tilde{X}_k :

$$\tilde{X}_k = i \Leftrightarrow X_{k\tau} \in A_i.$$

The process \tilde{X}_k describes the *snapshot dynamics* of X_t with lag time τ between the partition sets A_1, \dots, A_n . The so defined process is not necessarily Markovian, since $(G(\mathcal{P}^k))_k$ is in general **not** a semigroup.

see ..

Rebinding in a Double Well Potential

Let X_t be the Markov process corresponding to the double-well potential $V(x) = (x^2 - 1)^2$. We consider a full-partition of the state space into two sets A and B around the local minima of the energy landscape, as shown in figure 1.1. We are interested if the induced process \tilde{X}_k inherits the Markovianity of X_t or if it contains any memory effects.

For a small lag-time $\tau = 0.1$ we compute the probability of \tilde{X}_k to make a transition from B to A in one time-step. We compare it to the probability of the same transition with the **additional** information of having been in B one time-step before/earlier. If the process was Markovian, then this additional information about the past should make no difference and thus, both probabilities should be equal. We compute them in terms of the original process X_t by

$$\mathbb{P}_\mu[X_{(k+1)\tau} \in A \mid X_{k\tau} \in B] \text{ and } \mathbb{P}_\mu[X_{(k+1)\tau} \in A \mid X_{k\tau} \in B, X_{(k-1)\tau} \in A].$$

Using the density functions v_B and v_{BA} , we get

v densities?

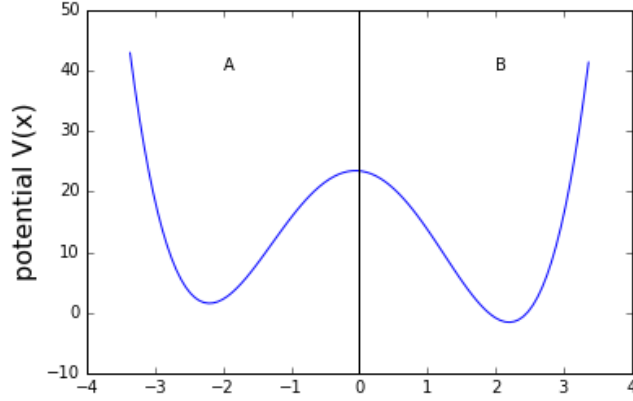


Figure 1.1: Full-partition of a double-well potential

$$\mathbb{P}_\mu[X_{(k+1)\tau} \in A \mid X_{k\tau} \in B] = \int_A v_B^\tau(x) dx = \dots, \quad (1.11)$$

$$\mathbb{P}_\mu[X_{(k+1)\tau} \in A \mid X_{k\tau} \in B, X_{(k-1)\tau} \in A] = \int_A v_{BA}^\tau(x) dx = \dots \quad (1.12)$$

So we see that for such a short lag-time τ , the process \tilde{X}_k is **not** independent of the past and hence **not** a Markov process. Equation (1.11) describes the probability to get from B to A , where “being in B ” could mean everything from “close to the transition region” to “far away from the transition region”. So this probability is averaged over **all** possible starting points in B . We compare it to (1.12), where being in A shortly before being in B increases the probability to return/recross to A again. This behaviour can be interpreted such that for a short time after a **transition**, the process is likely to be still inside of the **transition region**. In our example, the transition region is the area/region close to the maximum of potential energy. Thus, there is still an increased probability to return to the previous state.

This issue is called the *recrossing effect*, since additional memory leads to an increased probability to “recross” shortly after a transition. On the other hand, if we choose a large lag-time $\tau = 100$, then the past transition from A to B in (1.12) took place a long time ago. So we cannot certainly know if the process is still in the critical transition region; during that long lag-time it could also have been gone anywhere else. That means that the memory effect included in \tilde{X}_k becomes smaller for larger lag-times and thus can be considered as a *short-time memory*.

?

recross the barrier?

Comparison to Markov State Model

After having observed the *recrossing effect* as a memory effect when projecting the time-series(?) of a given continuous process onto a finite subspace, we want to compare that result to the corresponding *Markov State Model*.

def MSM

So far, we considered the process \tilde{X}_k belonging to the operator $G(\mathcal{P}^k)$. Now, let $(\hat{X}_k)_{k \in \mathbb{N}}$ be the Markov chain that is described by the transition matrix P_c , i.e. the matrix representation of the discretized transfer operator $G(\mathcal{P}) := G\mathcal{P}G$.

A desirable behaviour of this model would be that \hat{X}_k and \tilde{X}_k have the same trajectory when started on the same initial distributions \hat{X}_0 and \tilde{X}_0 . It will turn out that this is normally not the case. This question is visualized in diagram 1.2. Does it make a difference if we first project the process and then propagate it and vice versa?

$$G(\mathcal{P}(\tau)) \hat{=} \begin{array}{ccc} \mathcal{P}(\tau) & \xrightarrow{\tau \rightarrow \tau k} & (\mathcal{P}(\tau))^k \\ \downarrow \text{proj.} & & \downarrow \text{proj.} \\ P_c(\tau) & \xrightarrow{\tau \rightarrow \tau k} & (P_c(\tau))^k \end{array}$$

Figure 1.2: Projecting/propagating a transfer operator (non-commutative)

In general, this diagram is **not** commuting and hence, in general we have

$$(P_c)^k \neq (P^k)_c.$$

For the example of a full-partition discretization, we know from section 1.3 that the resulting Markov State Model is a Markov chain. Thus, we have a Markov chain \hat{X}_k as a model for the non-Markovian process \tilde{X}_k , so it is clear that there is a discretization/iteration error.

That is also where the term *Markov State Model* comes from. We are describing the non-Markovian process \tilde{X}_k by a Markov chain \hat{X}_k . Originally, processes have been clustered *set-based*, i.e. based on a full partition and thus always resulting in a Markov chain. In chapter 2, we will see that the *function-based* approach yields better results and hence is the current state of the art. Then, the Markov State Model is not necessarily Markovian, as we already know from $P_c = TS^{-1}$ in theorem 1.17.

Discretization Error (Density Propagating Error/Iteration Error)

We will describe here shortly how the discretization error can be estimated. For our purposes that will not play an important role, since later we will be able to perform a projection (with convenient membership functions) s.t. this error vanishes.

what about eigenv. err.?

membership fct = part. of unity for clustering

The maximal possible error between the distributions of \tilde{X}_k and \hat{X}_k after k time-steps is (independently of initial distribution) given by

which norm?

$$E(k) = \|G(\mathcal{P}^k) - (G(\mathcal{P}))^k\|.$$

Theorem 1.22. *Assume the discrete/dominant spectrum of a transfer operator \mathcal{P} is given/denoted/ordered by $1 = \lambda_0 > \lambda_1 \geq \dots \geq \lambda_n$. Then the projection error can be bounded from above in terms of the second-largest eigenvalue by*

$$\|(G(\mathcal{P}))^k - \Pi_0\| \leq \lambda_1^k,$$

where Π_0 is the orthogonal projection of

For a proof, see Schütte and Sarich[30, p.72]. In the next chapter we will see further/deeper relations between the spectrum of the transfer operator and the long-time behaviour of the process. We will see how to choose partition for a MSM s.t. the approximation error becomes small/vanishes.

def MSM, P_c
?

Conclusion

We have to distinguish between two kind of “errors” that can occur:

- Rebinding Events: Projection can include some kind of memory effect
- Iteration Error: Deviation of $G(\mathcal{P}^k)$ and $(G(\mathcal{P}))^k$

2 Dominant Structures

In section 1.3, we introduced a method to reduce the dimension of a Markov process by projecting it onto a smaller state space. But we don't know yet how to choose a partition of unity such that the Galerkin projection yields a reasonable Markov State Model, in the sense that important properties of the original process are inherited. The answer to the question, what an important property is, can differ from case to case. In this thesis, we are particularly interested in the long-time behaviour of a process, which is often influenced by some *dominant sets*. They can be characterized by the concept of *metastability*. We will see why it makes sense to project a process onto its metastable sets and, in order to detect them, analyze their relation to the dominant spectrum of the transfer operator. We will also see that the optimal metastable decomposition is not sharp/crisp but soft/fuzzy.

For a nonreversible process, such dominant structures can sometimes be given in terms of *dominant cycles* instead of dominant sets. We will introduce these structures as well and give a short outlook how they can be detected.

2.1 Metastability

There exist several different definitions of metastability. Shortly said, metastability is the property of a process that its state space consists of subsets/regions such that transitions between these subsets are rare events while the duration of stay inside of each of them is rather long. Some possible characterizations of that behaviour are based on large hitting times or small exit rates, see Schütte and Sarich[30, chapter 3], where a good overview of the most common definitions can be found.

Mathematical concept of metastability

In order to describe the concept of metastability, it is a good way to start with so called *stable* or *invariant subsets*. A measurable subset $A \subset E$ of the state space of a Markov process X_t is called stable or invariant if it cannot be left, i.e. if $\mathbb{P}(X_t \in A \mid X_0 \in A) = 1$ for all t . Analogously, we can define a *metastable* or *almost invariant subset* as a subset in which the process will stay for a very long time before exiting it into any other subset, that is $\mathbb{P}(X_{t_f} \in A \mid X_0 \in A) \approx 1$ for a convenient timescale t_f . Thus, a

full partition A_1, \dots, A_m of the state space E is called *metastable* if

$$\sum_{k=1}^m \mathbb{P}_\mu(X_{t_f} \in A_k \mid X_0 \in A_k) \approx m. \quad (2.1)$$

Then each of the sets A_k is almost invariant with respect to timescale t_f ; the probability to stay in one of the partition sets being started there is almost 1, while the probability to change between any two different partition sets is almost 0. Such a partition is also called a *metastable decomposition*.

Obviously, being “close to 1” or “close to m ” are rather vague statements. But that lack of concreteness will be eliminated later, since we will only be interested in the “best” metastable decomposition. That means that we want to obtain a decomposition where the sum (2.1) is as close as possible to m , or equivalently the probability to stay inside of a metastable set is as close as possible to 1. Also the choice of the timescale t_f is not specified in general and will depend on the particular system in consideration. Hence, the only parameter in (2.1) that has to be determined is the number m of subsets we are looking for.

joint
metasta-
bility?

Metastability in Molecular Dynamic Systems

Metastability is a very important concept for stochastic processes that describe the movement of atoms or molecules in space. Such processes have the characteristic behaviour to oscillate or fluctuate around equilibrium positions on the smallest time scales (about one femtosecond). In contrast to these fast oscillations, the process often stays inside of a certain region, also called *conformation*, for a long time before switching to another region (nano- or millisecond time scale). Conformational transitions are rare events. Thus, these conformations correspond to the above definition of metastable sets if we choose a convenient timescale.

Brownian
motion?

Figure 2.1 depicts/shows such a behaviour on the example of the dihedral angle of a molecule, taken from Weber[38]. This dihedral angle can take values between $+45^\circ$ and -45° . Thus, the process acts on a continuous state space and has infinitely many states. We see that there are two regions (highlighted red and blue) where the process stays for a rather long time and oscillates **inside**. Transition between these regions don’t occur often. Thus, these two conformations can be identified as metastable sets.

As transitions between metastable sets are rare events, we need to make long-time simulations of a process in order to get informations about these conformational changes. But long-time simulations of such large systems are not feasible in reasonable time even with the best computers nowadays, see Anton[31] or its successor Anton2[32].

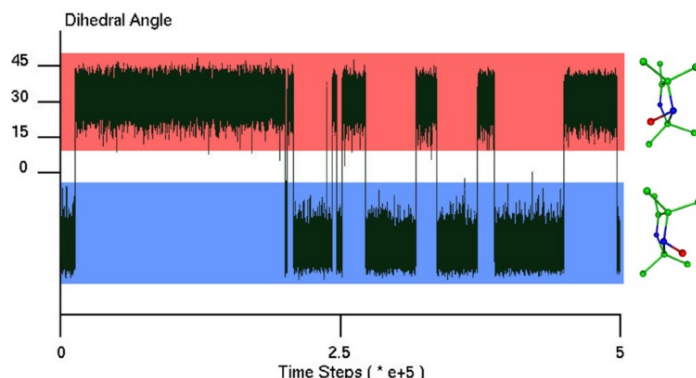


Figure 2.1: Example of a molecule with two (metastable) conformations

Hence, in order to be able to compute some long-time simulations of a given MD system, a reduction of complexity is needed. This can be achieved by a clustering of the state space via a Galerkin projection as depicted in section 1.3. Different states will be clustered appropriately such that we get a process on a smaller state space.

This MD point of view also motivates the following terminology. A state in the original state space is called a *micro state*, as it is a state considered on the microscopic or atomistic level. In order to get a smaller state space, micro states are grouped together and such a cluster is called a *macro state*, since we are now considering the process on a macroscopic level (cannot distinguish between smaller states/atoms anymore).

For instance, the spatial coordinates of a single atom could be considered as a micro state, while the corresponding macro state is a cluster of several atoms. If we are working on this smaller state space, we cannot distinguish anymore between the single atoms of the cluster (forget information).

Clustering into metastable sets

The question **how** to cluster the process (i.e. how to choose the partition of unity for Galerkin projection) such that it maintains the long-time behaviour of the original process can be answered easily with the following intuitive approach: As the long-time behaviour of a process is described by the metastability of the system, we choose the metastable sets as clustering sets. More clearly, we create a new process where each macro state corresponds to one of the metastable sets. The transition probabilities of the clustered process/reduced model should correspond to the transition rates of the original process between its metastable sets in order to represent the correct (long-time)

behaviour of the process.

As metastability is determined on long timescales, the reduced/projected process should maintain the long-time behaviour of the original process, but *forget* about its short-time transitions, i.e. transitions inside of a conformation/metastable set.

Since there is not one unique metastable decomposition of the state space, we need to find a decomposition which is in some sense “the best”; then we can use it to create a reduced model. In sections 2.2 and 2.3 we will see how to find such a decomposition.

Most importantly, the projected/clustered process will have the desired property of a reduced dimension/complexity since the model acts on a smaller state space while maintaining the crucial property of the original process (transitions between metastable sets = long-time behaviour). So the computation effort for (long-time) simulations is definitely decreased. Furthermore, we get a better overview of the system, since it is always easier to consider a process on a few states in comparison to a process on a very large or even continuous state space. Since fast/short-time transitions (transitions inside conformation/metastable set) are not our research goal, we just omit these (at least for our case) superfluous informations. But there is also a disadvantage, as already mentioned in section 1.4, by projecting a process it can lose its Markov property.

+ proj. error
possible

2.2 Spectral Approach

In this section we will see that the spectrum of the transfer operator is highly connected to the metastability of the corresponding Markov process. Namely, the number of metastable sets can be determined by the number of eigenvalues close to 1 and the corresponding eigenfunctions induce a metastable (full) decomposition.

Existence of dominant eigenvalues

We consider the transfer operator $\mathcal{P} := \mathcal{P}^t$ of a Markov process for some fixed t in the Hilbert space $L^2(\mu)$. We are interested in *dominant eigenvalues* of \mathcal{P} , that is large eigenvalues which are close to 1 and separated from the rest of the spectrum. The *discrete spectrum* $\sigma_{\text{discr}}(\mathcal{P})$ is the set consisting of all eigenvalues $\lambda \in \sigma(\mathcal{P})$ that are isolated and of finite multiplicity. The *essential spectral radius* $r_{\text{ess}}(\mathcal{P})$ is defined as

$L^2, L^1?$

$$r_{\text{ess}}(\mathcal{P}) = \inf\{r \geq 0 \mid \lambda \in \sigma(\mathcal{P}) \text{ with } |\lambda| > r \text{ implies } \lambda \in \sigma_{\text{discr}}(\mathcal{P})\}.$$

The existence of dominant eigenvalues requires that the continuous part of the spectrum is bounded away from the dominant elements of the discrete spectrum. To ensure that the process we are considering actually possesses metastable sets, we need to pose some conditions on the spectrum of the transfer operator:

C1 The essential spectral radius of \mathcal{P} is less than one; i.e. $r_{\text{ess}} < 1$.

C2 The eigenvalue $\lambda = 1$ of \mathcal{P} is simple and dominant; i.e. $\eta \in \sigma(\mathcal{P})$ with $|\eta| = 1$ implies $\eta = 1$.

We will not go into further details for which processes the two above conditions are fulfilled; some criteria for it can be found in Huisinga[14, chapter 4]. Since these conditions are required for the later investigations, we will just assume that they are true.

We need condition **C1** to ensure that the continuous part of the spectrum is bounded away from the discrete eigenvalues. Otherwise they would not be dominant anymore and the process would be rather rapidly mixing than having any metastable sets/dominant structures. Condition **C2** however is important because the state space of a transfer operator with more than one eigenvalue of absolute value 1 can be decomposed into invariant sets, that is subsets which cannot be left. But that case is not interesting for us. Instead, we want to know more about **almost** invariant sets and their critical transition regions.

C2 = ergodic?

Huisinga?

Theorem 2.1. (Schütte[30, theorem 4.16])

The transfer operator $\mathcal{P} : L^2(\mu) \rightarrow L^2(\mu)$ of a reversible process with properties **C1** and **C2** has the following spectrum:

$$\sigma(\mathcal{P}) \subset [a, b] \cup \{\lambda_n\} \cup \dots \cup \{\lambda_2\} \cup \{1\}$$

with $-1 < a \leq b < \lambda_n \leq \dots \leq \lambda_1 = 1$ and isolated, not necessarily simple eigenvalues of finite multiplicity that are counted according to multiplicity.

This theorem assures us the existence of a discrete set of dominant eigenvalues. In the following we will see that this property results in metastability.

Relation of dominant eigenvalues to metastable sets/Optimal Decomposition

Theorem 2.2. (Schütte[30, theorem 4.16])

The metastability of an arbitrary decomposition $\mathcal{D} = \{A_1, \dots, A_m\}$ of the state space \mathbb{X} can be bounded from above by

$$p(A_1, A_1) + \dots + p(A_m, A_m) \leq 1 + \lambda_2 + \dots + \lambda_m.$$

Theorem 2.2 provokes the question if there exists an *optimal* decomposition with highest possible metastability.

maybe ill-conditioned

The number of metastable sets is determined by the number of dominant eigenvalues.

This theorem shows/reveals the relation/connection between metastable sets and the eigenfunctions and eigenvalues of the transfer operator.

Relation of dominant eigenfunctions to metastable decomposition

Theorem 2.3. (Schütte[30])

Each single eigenfunction induces a metastable decomposition

Proof.

□

The zeros of an eigenfunction (respectively change of sign) induce a metastable decomposition of the state space. Different eigenfunctions result in different decompositions.

This theorem gives us a first demonstrative relation of the metastability of a process to its eigenvectors. This result is not yet optimal/very good. In the next sections, we will see that linear combinations of eigenvectors result in much better metastability.

eigenfcts vs
committor
functions

In figure 2.2, we get a good overview of that relation. We have a potential/energy landscape which has 4 energy minima, i.e. 4 regions where the process could be *trapped* such that it is hard to get outside again. The transition matrix shows this metastable behaviour since we can see 4 regions which large probabilities to stay inside and very small probabilities to go to a different region. Furthermore we see that the process has 4 dominant eigenvalues. Three of them have a change of sign, which induces the metastable decomposition.

Disadvantages

So the spectral approach is suitable/convenient to characterize metastability of Markov processes. But there are two disadvantages. 1: the result is only applicable on reversible processes, because real eigenvalues are only guaranteed if the transfer operator is self-adjoint. 2: eigenvector problem of the transfer operator has only global solutions.

see

Most of all, the previous approach in computing a metastable decomposition doesn't include/consider the transition regions of the process, so we need to refine/improve it.

Remark: the previous theorems and this entire section have been presented only to **see** the relation of the spectrum to the metastability of a system (which is certainly a strong relation!). This relation can be seen best for a full partition decomposition/clustering. But in the next chapter, we will get a clustering wrt metastability which yields even better results; it will still beinhalten the eigenvalues and eigenfunctions of the transfer operator. But it will be FUZZY!

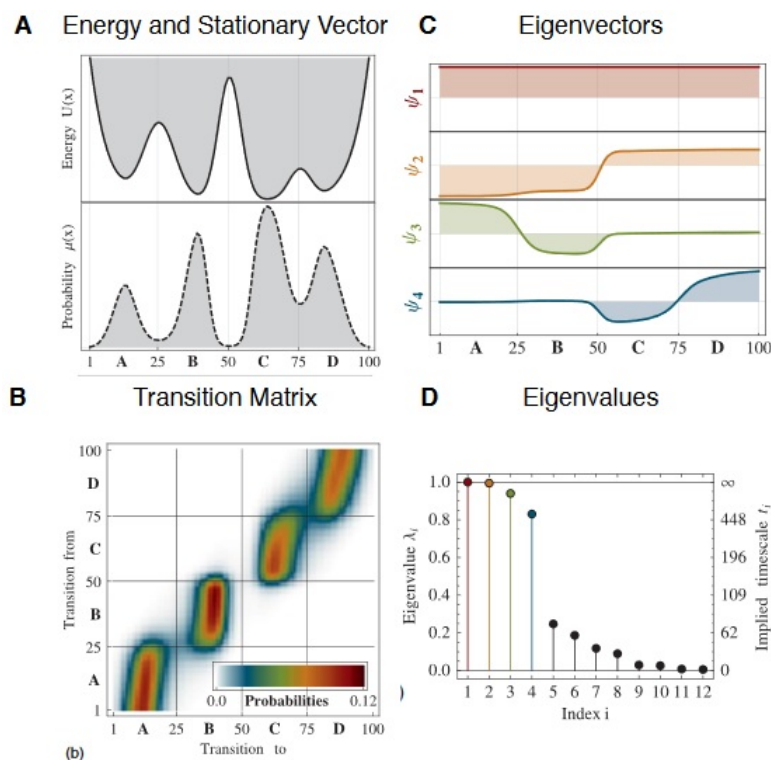


Figure 2.2: Relation of eigenvalues and eigenvectors to metastability of a process

2.3 Fuzzy Clustering

The above considerations result in a metastable full decomposition of the state space, that is each state is assigned to exactly one of the partition sets. We will see now, that there exist better solutions, considering/including the fact that transition regions can belong to several metastable conformations. So for this slightly more general approach, there may be some overlap in the assignment of states to metastable sets.

Set-based vs. Function-based Approach

An intuitive approach to decompose the state space would be to determine a certain number of metastable sets which form a full partition of the state space, such that each state is assigned to exactly one of the metastable sets. The problem with that approach is that also the transition regions of the process would have to be assigned to one of these partition sets. But why would you assign a state in a transition region to one adjacent metastable set and not to the/one other? So such an assignment is not a rigorous description of actual behaviour of the process.

Therefore, this *set-based* or *crisp approach* of decomposing the process has been replaced by the *function-based* or *fuzzy/soft approach*. That means that each state of the process is assigned with a certain “degree of membership” (\approx probability) to each metastable conformation. That is reasonable, because a state in a transition region (local minimum) cannot be assigned to a single conformation. Thus, these clustering-functions should be “overlapping”.

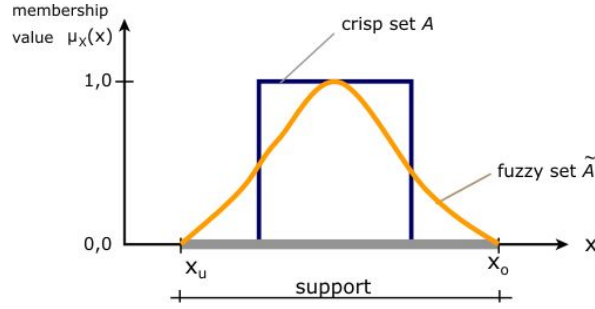


Figure 2.3: Crisp vs Fuzzy Sets/ Clustering

Fuzzy sets are sets whose elements have degrees of membership.

Membership functions

Assuming we have already determined that our process consists of n metastable sets (by knowing its n dominant eigenvalues). We will follow the approach of Weber[37] to define macro states as *overlapping partial densities*. They can be identified by membership functions $\chi_1, \dots, \chi_n : \mathbb{X} \rightarrow [0, 1]$. Each state of the original state space shall be assigned to the different macro states with a certain *degree of membership*.

Definition 2.4. (Membership Function [37])

The functions $\chi_1, \dots, \chi_n : \mathbb{X} \rightarrow [0, 1]$ are called *membership functions* if they fulfill

- $\chi_j(i) \geq 0 \ \forall i \in \mathbb{X} \text{ and } \forall j \in \{1, \dots, n\}$ (positivity),
- $\sum_{j=1}^n \chi_j(i) = 1 \ \forall i \in \mathbb{X}$ (partition of unity).

We can interpret membership functions as assigning each state i to a cluster/conformation j with a certain “probability” $\chi_j(i)$. For each conformation $j \in \{1, \dots, n\}$, there is a membership function χ_j , which determines the portion of the partial density w.r.t. the total density function. These membership functions form a partition of unity, ? in order to sum up to the total density.

The previous example of a full-partition discretization corresponds to the choice of characteristic functions $\{\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_n}\}$ as membership functions. They are also called *crisp* or *hard* membership functions, whereas the general (overlapping) membership functions are denoted as *fuzzy* or *soft*. The *crispness* of the clustering(?) can be measured by the matrix S . Nonoverlapping membership functions (char. fcts.) yield an overlap matrix $S = D^{-1}\langle\chi, \chi\rangle$ equal to the unit matrix. Overlapping membership functions yield a matrix with non-zero outer diagonal elements.

Membership functions can be used to decompose the space into metastable sets/clusters. In this case, the assignment of a state to a metastable set must not be unique, but a state can belong to different metastable sets with certain degrees, which can be interpreted as kind of probabilities. That model takes into consideration the existence of transitions regions which cannot be uniquely assigned to one macro state/metastable set. ?

Since membership functions form a partition of unity, we can apply the Galerkin projection as defined in section 1.3. If we choose the membership functions w.r.t. metastability, then we get a Markov State Model where each (macro) state is a metastable set of the original process.

so far: membership fct indep. of metast. set ?

Individual eigenfunctions \mathcal{X} do not overlap since they are orthogonal. But the membership functions χ_j as linear combinations of the dominant eigenfunctions might have an overlap.

$$P_c \chi_j \approx \chi_j \quad \forall j = 1, \dots, n.$$

$$T \chi_j \approx S \chi_j \quad \forall j = 1, \dots, n.$$

From the definition, membership functions like $\chi_1 = \chi_2 = 0.5$ are possible, but not interesting. Instead, they are often chosen to be **close** to a characteristic function, like shown in figure 2.3. That is reasonable, since it puts the emphasis of a conformation onto a certain region (high degree of membership) and maybe some adjacent parts (low degree of membership). For this reason, they are also often called *almost characteristic functions* in literature.

In section 1.3, we defined the Galerkin projection on an arbitrary partition of unity $\{\chi_1, \dots, \chi_n\}$, instead of defining it just on full partitions. Such a partition of unity (membership fcts.) induces a projected transfer operator P_c . Then the trace of P_c is referred to as *metastability* of the conformations $\{\chi_1, \dots, \chi_n\}$. This definition of metastability corresponds to the full-partition formulation of metastability (2.1), since the projected process has the state space $\{\chi_1, \dots, \chi_n\}$. Therefore, the diagonal of P_c consists of the *holding probabilities* of the conformations.

membership fct, cluster, metastability, macro state

Statistical Weights

For each macrostate we can define a statistical weight

$$w_i = \langle \chi_i, \mathbb{1} \rangle_\mu = \int_{\mathcal{X}} \chi_i(q) d\mu(q),$$

which describes the *portion* of a membership function to the total density function. $D = \text{diag}(w_1, \dots, w_n)$ is the diagonal matrix of the statistical weights of the membership functions. Then $T = D^{-1} \langle \chi, \mathcal{P}(\tau) \chi \rangle_\mu$, compare theorem 1.17.

what for?

?

Perron Cluster Analysis

The term *Perron Cluster Analysis* denotes the objective of clustering a Markov process into metastable sets using the *Perron eigenvalues* respective *Perron eigenfunctions*, which means eigenvalues close to 1 and the corresponding eigenfunctions. Perron Cluster Analysis respectively its algorithmic implementation PCCA (*Perron Cluster Cluster Analysis*) has been developed by Deuffhard et al[8] which used the sign structure of the dominant eigenvalues of the transition matrix. That approach has been improved by Deuffhard and Weber[9] who transformed the system of eigenvectors into a system of membership functions which results in a soft/fuzzy clustering of the state space of the original process; their algorithm is called PCCA+ (*Robust Perron Cluster Analysis*). Originally, PCCA+ was formulated only for discrete Markov chains, but Weber[38] extended it even on continuous processes.

finite/cont.
vec./fct.

set-based ap-
proach?

We consider the set of dominant eigenvalues $\{\lambda_1, \dots, \lambda_n\}$ with the corresponding set of eigenfunctions $\mathcal{X} = \{\chi_1, \dots, \chi_n\}$. They fulfill the eigenvalue problem $\mathcal{P}(\tau)\mathcal{X} = \mathcal{X}\Lambda$ of the transfer operator $\mathcal{P}(\tau)$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. The set of membership functions $\chi = \{\chi_1, \dots, \chi_n\}$ can be built as a linear combination $\mathcal{X}\mathcal{A}$ of the dominant eigenfunctions, that is

$$\chi_j(q) = \sum_{i=1}^n \mathcal{A}_{ij} \chi_i(q), \quad j = 1, \dots, n. \quad (2.2)$$

Here, $\mathcal{A} = \{\mathcal{A}_{ij}\}_{i,j=1,\dots,n} \in \mathbb{R}^{n \times n}$ is a real matrix which should be chosen in such a way that the resulting membership functions fulfill the required properties/constraints; i.e. positivity and partition of unity. There are infinitely many transformations \mathcal{A} of the eigenvectors resulting in a soft membership matrix χ satisfying the positivity and partition of unity constraints. Consequently, we have to determine the transformation \mathcal{A} that satisfies some optimality condition. The algorithm PCCA+ computes the matrix \mathcal{A} as the solution of a convex maximization problem, see Weber[37].

optimal
metastabil-
ity?

only conv.
lin.comb.?

Weber shows in [38] that for choice $\chi = \mathcal{X}\mathcal{A}$ (any linear combination of the eigenfunctions \mathcal{X}), the discretization error of the Galerkin projection vanishes and hence, diagram 1.2 commutes. In particular, such membership functions preserve the Markov property of the process.

If we project a finite process, then the above formulation results in a *membership vector matrix* χ , which is the result of a linear combination of the *eigenvector matrix*.

Measuring/Maximizing crispness of

The *crispness* of .. can be measured by the matrix S .

The metastability of a process is defined via the trace of P_c , but can also be measured via its determinant. If the metastability is high, then $\det(P_c)$ is close to 1. More clearly, see ... by multiplication of the determinant we have

$$\det(P_c) = \det(S) \det(\Lambda).$$

Thus, in order to increase the metastability of a system, both determinants need to be high. $\det(S)$ is maximized if the linear combination $\chi = \mathcal{X}\mathcal{A}$ is as *crisp* as possible.

2.4 Dominant Cycles

When it comes to computing a metastable decomposition for a nonreversible process, we have to face the problem that the eigenvalues/eigenvectors might be complex valued and thus, PCCA+ is not applicable. One possibility/alternative/way to circumvent this problem is to consider the real Schur decomposition of the matrix(?) instead of its spectral decomposition. Then, we can apply PCCA+ to the real Schur vectors of the matrix instead to its eigenvectors. This approach is feasible/possible, since the real Schur vectors span the same subspace as the corresponding complex Schur vectors and those span the same subspace as the corresponding eigenvectors. of transf.op.? see ...

Furthermore, a nonreversible process can contain other dominant structures than just *metastable sets/conformations*. It can as well contain *metastable cycles*, that is subsets with a cyclic behaviour and a high probability to stay inside of such a cycle for a long time.

NESS processes

Definition 2.5. (NESS process)

A Markov process is called *nonequilibrium steady state (NESS)* process if it is nonreversible, but still has a steady state, given by an invariant measure μ w.r.t. which the process is ergodic. what is that?

$$p_\tau(A, B) = p(\tau, A, B) = \mathbb{P}(X_\tau = B \mid X_0 = A)$$

As a NESS process is nonreversible, there are regions where the detailed balance equation is not fulfilled, i.e. there is an effective probability flow $p(\tau, A, B) - p(\tau, B, A) \neq 0$ between some subsets $A, B \subset S$ of the state space.

Flow of a process

In the following we consider an irreducible and aperiodic (i.e. ergodic) Markov chain on the finite state space $S = \{1, \dots, n\}$ given by the transition matrix P . Since this Markov chain is irreducible and aperiodic, it possesses a unique invariant measure μ that is positive everywhere. Then μ is the normalized eigenvector of P for the unique eigenvalue $\lambda = 1$.

why only finite?

see ..

Definition 2.6. (Flow Matrix)

The probability flow associated to a Markov process is given by the flow matrix

$$F = DP,$$

where P is the transition matrix of the process and D the diagonal matrix $D_{ii} = \mu_i$ with the entries of the invariant measure μ .

So the (steady state) probability flow from state i to j is given by $F_{ij} = \mu_i P_{ij}$. If the process is reversible, the flow matrix F is symmetric due to the detailed balance equation. For a NESS process, F is not symmetric since there are states $i, j \in S$ with $F_{ij} \neq F_{ji}$.

Cycle Decomposition

This flow must be decomposable into elementary cycles.

why?

Definition 2.7. (Cycle of a process)

A k -cycle γ on S is an ordered sequence (up to cyclic permutations) of k connected states $\gamma = (i_1, \dots, i_k)$ with length $|\gamma| = k$, i.e. the probability to get to the next state is always positive: $\mathbb{P}_{i_j, i_{j+1}} > 0$ and $\mathbb{P}_{i_k, i_1} > 0$. Cycles without repetition/self-intersections are called simple cycles. The set of all simple cycles is denoted by \mathcal{C} .

1-step prob.?

We want to make a cycle decomposition of the flow F , see Kalpazidou[15].

Definition 2.8. (Cycle/flow Decomposition)

A collection $\mathcal{C}_+ \subset \mathcal{C}$ of cycles γ with real positive weights $w(\gamma)$ is a flow decomposition if for every edge $(i, j) \in S^2$ we have

$$F_{ij} = \sum_{\gamma \supset (i,j)} w(\gamma),$$

where $(i, j) \in \gamma$ if the edge (i, j) is in γ .

In order to make sense in a probabilistic context, we define the weight w of a cycle γ in the following way. Given a (realization of?) Markov chain $(X_t)_{t \in \mathbb{T}}$, we count the number of times N_T^γ the process passes through a cycle γ up to time T .

Definition 2.9. (Weight of a Cycle)

$$w(\gamma) = \lim_{T \rightarrow \infty} \frac{N_T^\gamma}{T}.$$

Since we are considering/assuming an ergodic process, this limit exists a.s.

jian qian

Dominant cycles/sets

We will see that dominant cycles have similar properties as dominant sets for reversible processes, i.e. large eigenvalues with $|\lambda| \approx 1$. But now the eigenvalues are lying in the complex plane and might be non-real (pairs of complex eigenv.).

Definition 2.10. (Dominant Cycle)

= metastable cycle?

So a cycle is dominant if there is a high probability inside the Markov chain to follow this cycle.

Dominant structures will be defined utilizing the dominant Schur vectors of the transition matrix instead of its eigenvectors. A membership matrix can be defined as a linear combination of these leading Schur vectors (spectral clustering with PCCA+).

Schur Decomposition

We have the same situation/aim as in the previous sections: we have a Markov process on a large state space S and we want to decompose it into a smaller state space consisting of clusters that belong to dominant structures of the process.

Definition 2.11. (Schur Decomposition)

Let $P \in \mathbb{R}^{n \times n}$ be a transition matrix. Then it can be written as

$$P = XRX^{-1},$$

where X is a unitary matrix and U is an upper triangular matrix, which is called a *Schur form/matrix/decomposition* of P .

Since R is similar to P , both matrices have the same spectrum. Since R is triangular, their eigenvalues are the diagonal entries of R .

A Schur Decomposition is not unique. As P is a real matrix, its non-real eigenvalues ...
 come in complex conjugate pairs. This fact can be used to build a *real Schur form*,
 where X and R are both real matrices. But then R is no longer triangular, but only
quasi-triangular, allowing 2×2 -blocks on its diagonal. The eigenvalues of the 2×2 -blocks picture
 are exactly the complex conjugate eigenpairs of P . see ..

Theorem 2.12. (Real Schur Decomposition [33, Exc. I.3.24])

If $A \in \mathbb{R}^{N \times N}$, then there exists an orthogonal matrix $U \in \mathbb{R}^{N \times N}$ such that

$$U^T A U = T,$$

where T is block-triangular with 1×1 and 2×2 -blocks on its diagonal. The 1×1 -blocks
 contain the real eigenvalues of A and the eigenvalues of the 2×2 -blocks are the complex
 eigenvalues of A .

Definition 2.13. (Schur Vector)

Let $\tilde{R} \in \mathbb{R}^{m \times m}$ be a submatrix of R (top left part of R). Then

$$P = \tilde{X} \tilde{R} \tilde{X}^{-1},$$

where $\tilde{X} \in \mathbb{R}^{n \times m}$ consists of the first m columns of X . These vectors will be denoted
 as the dominant Schur Vectors of P . Schur Values? why?

Using PCCA+

Djurdjevac Conrad et al[10] propose an algorithm in order to determine/get the desired
 membership vectors χ .

- i) Compute a real Schur decomposition (\tilde{X}, R) of ...
- ii) Sort the Schur values and the 2×2 -blocks such that they are in a descending
 order
- iii) Determine the submatrix \hat{X} and solve PCCA+ equation (...) in order to get the
 membership functions χ

Computation of metastable cycles/sets

Fackeldey and Weber[11]

3 Rebinding Effect in a Given Kinetics

In this chapter, we will examine a particular type of molecular systems, namely receptor-ligand systems, consisting of special molecules, so called receptors and ligands. These molecules can interact in such a way that under certain conditions they can *bind* to each other and afterwards *dissociate* again. Such a system is originally Markovian and can be described by a transfer operator. However, by projecting this operator onto a finite-dimensional state space, the Markov property of the process may be spoiled.

We explain this memory effect, the so called *rebinding effect*, and examine its influence on the quantity of binding events of a receptor-ligand-system. We show how this effect can be measured with tools we know from chapter 1. We are particularly interested in analyzing an optimization problem to find a lower bound for the rebinding effect.

This chapter highly complies with Weber and Fackeldey[40]. Additionally, we introduce some fundamental definitions and notations from biochemistry, which are necessary for the understanding of receptor-ligand-systems. The importance of such systems becomes clear by a description of its application for drug design.

3.1 Receptor-Ligand System

We present a special molecular system, the *receptor-ligand-system*, and model it mathematically using a differential equation. We discuss the so called *rebinding effect*, a memory effect occurring in this model, and set it in relation to the recrossing effect known from section 1.4. We explain its relevance in the application of *drug design*.

Molecular Dynamics vs Molecular Kinetics

A *molecular system* consists of atoms that are connected by *covalent bonds*. bla bla. The potential energy function, or *energy landscape*, of a molecular system results in a dynamical behaviour on different timescales. The fastest timescales (vibrations of covalent bonds) are around 10^{-15} seconds. bla. up to nanoseconds or, for protein folding, up to microseconds or seconds or even longer.. see ..

Receptors and Ligands

references

In biochemistry, a *receptor* is a molecule, often a protein, that is usually located on the surface of a cell and can receive signals from outside the cell. A molecule that has the ability to *bind* or *associate* to a receptor is called a *ligand*. Each receptor will only bind with ligands of a particular structure, which is often referred to as the “key-lock principle”. Both receptor and ligand need to have specific complementary geometric shapes that fit exactly into one another, as exemplarily depicted in figure 3.1.

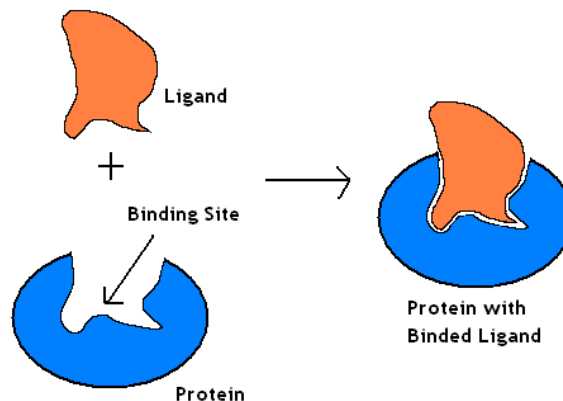


Figure 3.1: Ligand (“key”) binds to a receptor (“lock”). Their shapes fit together.

binding site

Such a binding between a receptor and a ligand can *activate* (“unlock”) the receptor by producing some kind of a chemical signal and thereby provoke a physiological response. For instance, that could be a conformational change in a protein, caused by a hormone binding to it. However, instead of engaging in the actual physiological consequences of a binding, we focus on the **act** of binding events.

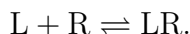
?

The action of binding is typically reversible¹ through *dissociation* of the involved receptor and ligand. Ligand binding is a *chemical equilibrium* process, which means that the reaction rates of the binding and dissociating events are equal, once this equilibrium is reached. From then on, the concentrations of the reactants (ligands) and the products (complexes) are constant. It is a *dynamic equilibrium*, since reactions take place, even though no net change in the concentrations can be observed.

The binding behaviour of a receptor-ligand system is formalized as follows. Ligands (L) can bind to receptors (R) and form receptor-ligand complexes (LR) which can

¹We remark that in this context, *reversible* means that a ligand can bind and unbind to a receptor, that is the reaction can run forward and backward. In contrary to the mathematical reversible, which means that a process behaves **equally** when running backwards in time.

dissociate again into its original components. These reactions can be represented by



Being a process in chemical equilibrium, the law of mass action states that the ratio between the concentration of reactants and products is constant. The corresponding *dissociation constant* k_d is given by

$$k_d = \frac{[L] \cdot [R]}{[LR]},$$

where $[L]$ is the concentration of free/unbound ligands, $[R]$ is the concentration of unoccupied receptors and $[LR]$ is the concentration of receptor-ligand complexes. This constant is used to describe the *binding affinity* between a ligand and a receptor, that is how strongly/tightly the ligand can bind to his particular receptor. If the dissociation constant is small, then there are relatively many complexes in comparison to unbound molecules, and for this reason, the binding affinity between the ligand and the receptor is high. The *association constant* k_a is just the inverse of the dissociation constant

$$k_a = \frac{[LR]}{[L] \cdot [R]}.$$

There are different factors which can influence the binding affinity of a process. It depends on the nature of the constituent molecules, like their shape, size and possible charge. The binding affinity of a particular ligand-protein interaction can also change significantly with solution conditions (e.g., temperature, pH and salt concentration).

In general, high-affinity binding results in a higher degree of occupancy for the ligand at its receptor binding site than is the case for low-affinity binding; the residence time (lifetime of the receptor-ligand complex) does not correlate. ??

Mathematical Model of Receptor-Ligand-System

Starting from the reaction equation (3.1), we can deduce that the ligand can be found in two different (macro) states: “unbound” (L) or “bound” (LR). Then the probabilities of the ligand to be in one of these states can be described by the probability vector $x^T = \frac{1}{s}([L], [LR])$, where $s = [L] + [LR] = \text{const.}$ is the normalization constant. This leads to an ordinary differential equation

$$\dot{x}^T = x^T Q_c.$$

The matrix Q_c consists of the rates of reaction,

$$Q_c = \begin{pmatrix} -k_a[R] & k_d[LR] \\ k_a[R] & -k_d[LR] \end{pmatrix},$$

where k_a and k_d are the association and dissociation constants. It corresponds to the transition rate matrix of a Markov chain, that means it describes a **memoryless** process. We will later see that this mathematical description of a receptor-ligand-system is not accurate, since in fact, such a process **will** have some kind of memory.



Figure 3.2: Two possible states (unbound, bound) of a ligand-binding-system.

The two possible states for the easiest case of a ligand-binding-system consisting of one receptor and one ligand are depicted in figure 3.2. We notice that the spatial arrangement of the receptor and the ligand in the unbound case is **not** included in the above model. Therefore, we cannot distinguish if, at a given time, the receptor and the ligand are close to each other or not.

Rebinding Effect

In fact, a stochastic process describing a receptor-ligand system is **not** necessarily Markovian. That is due to the spatial arrangement of the system after the dissociating of a receptor-ligand-complex took place. Shortly after such a dissociating, it is more likely that the corresponding receptor and ligand will bind again, since they are still close to each other. Such a binding shortly after being dissociated is called a **rebinding**. The memory effect which thereby occurs is called **rebinding effect**. On large timescales, this effect will vanish since the favorable spatial situation is not necessarily given anymore and the system will be rather mixed again. Thus, Markovianity can be spoiled by the rebinding effect. It is depicted in figure 3.3.

reb. eff. :=
mem. eff.

The rebinding effect and its occurrence in natural science has been described and analyzed by several authors[12, 36]. In chemistry, it has been discussed in the context of clustered receptors and clustered ligands[5]. Recently, there has been efforts to describe the rebinding effect mathematically, see Weber et al[39, 40].

multivalence

The rebinding effect has been discussed to increase the binding affinity of a process.

?

Rebinding Effect vs Recrossing Effect

The characterization of the rebinding effect reminds us of the recrossing effect, as described in section 1.4. There, we considered the projection of a process onto a

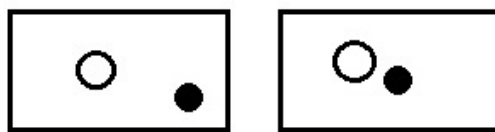


Figure 3.3: Rebinding Effect. Two configurations of a system consisting of one receptor (white) and one ligand (black), represented by the same state (“unbound”). Left: an arbitrary spatial arrangement of the receptor and the ligand in the unbound state. Right: spatial constellation of receptor and ligand shortly after their dissociation. The ligand is still **close** to the receptor, increasing the probability of a fast **rebinding**.

finite state space. This projected process (MSM) was described by a transition matrix, thus being a Markov chain, even though the process actually contained a (short-time) memory. The same phenomena occurs with the rebinding effect. We have a process which is modelled by a Markov chain, even though the process actually has a memory. We can interpret the two states of the ligand-binding-system as macro states resulting from the projection of a process on a larger state space (for instance, including more informations about the spatial situation of the receptor and the ligand). In this case, the rebinding effect originates from the loss of information caused by the projection and therefore, qualitatively corresponds to the recrossing effect.

it. error vs
reb. eff.

quant.
same?

Thus, the rebinding effect is just a special case of the recrossing effect. The different nomenclatures are justified by the use in their original context. While the term recrossing effect denotes the act of **recrossing** a energy barrier, the term rebinding effect denotes the **rebinding** of a receptor-ligand-system. Such a system can also be represented by a potential energy function, having energy minima around the bound (closest possible distance) and unbound state (farthest possible distance).

In chapter 2, we learned that a crisp clustering leads not to the best result and that we should chose a fuzzy clustering instead. Thus, we are going to include the spatial situation of the system by introducing degrees of membership, which could be interpreted as intermediate state such as “almost bound”. In this case, an unbound state with a high degree of membership to the bound state could correspond to a ligand which is close to a receptor, e.g. shortly after dissociating.

In the next sections, we are going to quantify this effect by embedding the whole molecular system into the mathematical framework established in the first two chapters.

Application: Drug Design

The term *drug design* denotes the development of new medications based on the knowledge of a biological target, playing the role of the receptor. Drug design is basically about designing a molecule which is complementary in shape and charge to the biomolecular target and therefore will bind to it, see Strømgaard et al[34]. More precisely, drug design describes the design of ligands, that is molecules that will bind tightly to the given target, see Tollenaere[35]. In general, we can distinguish between following two most common functionalities of drugs.

structure-
based vs.
rational DD

- **Activators** are able to activate, or even deactivate, a receptor and result in a strong biological response. An example for such a drug is morphine, which acts directly to the central nervous system, mimics the actions of endorphins and thereby reduces pain.
- **Inhibitors** bind to a receptor without activating it. Though, as they “block” the binding sites of receptors, they prevent possibly disease causing ... objects .. to bind. A well-known example are protease inhibitors, a class of antiviral drugs that are widely used to treat HIV and hepatitis C.

Independently of the fact whether a drug activates or inhibits receptors, a high binding affinity is required in order to be an efficient drug. The central dogma of receptor pharmacology (“occupation theory”) is that a drug effect is directly proportional to the number of receptors that are occupied. Furthermore, a drug effect ceases as a drug-receptor complex dissociates. Thus, a low binding affinity needs to be compensated by a higher concentration of ligands. Though, high concentrations should be avoided, because of possible side effects. Accordingly, the most fundamental goal in drug design is to predict whether a given molecule will bind to a target and if so how strongly.

Multivalence

covalent?

Bivalent ligands consists of two molecules connected by an (inert?) linker.

In general, multivalent ligands have a higher binding affinity because of the favorable spatial arrangement. If one of the ligands dissolves, then the other (connected) ligands still hold it in place (close to receptor).

3.2 Molecular Kinetics as a Projection

In this section, we will basically embed the mathematical concepts/results of chapter 1 into a chemical/physical context in order to get a rigorous description of molecular (dynamic/kinetic?) systems. When considering such systems, we can distinguish between two point of views: we will see how we can get from the *microscopic* or *atomistic* to a *macroscopic* scale/point of view by a projection.

Micro States

A micro state of a molecular system with N atoms can be represented in a $6N$ -dimensional *phase space* $\Gamma = \Omega \times \mathbb{R}^{3N}$, consisting of the *configurational space* $\Omega = \mathbb{R}^{3N}$ and the *momentum space* \mathbb{R}^{3N} . In the following, we consider systems in *thermodynamical equilibrium*. One possible model is given by the *Boltzmann distribution* $\pi : \Omega \times \mathbb{R}^{3N} \rightarrow \mathbb{R}$, a probability distribution assigning to each micro state a probability depending on its energy and temperature, see McQuarrie[19]. It can be expressed as

$$\pi(q, p) = \frac{1}{Z} \exp(-\beta H(q, p)), \quad (3.1)$$

where $\beta = 1/(k_B T)$ is the inverse of the temperature T multiplied with the Boltzmann constant k_B and $Z = \int_{\Gamma} \exp(-\beta H(q, p)) \, d(q, p)$ is the normalization factor. The Hamilton function denoted by H is given by $H(q, p) = K(p) + V(q)$, the sum of the kinetic energy $K(p)$ and the potential energy $V(q)$. Thus, the Boltzmann distribution π can be decomposed into $\pi = \pi_p \pi_q$,

$$\pi(q, p) = \underbrace{\frac{1}{Z_p} \exp(-\beta K(p))}_{\pi_p} \cdot \underbrace{\frac{1}{Z_q} \exp(-\beta V(q))}_{\pi_q},$$

where $\pi_p : \mathbb{R}^{3N} \rightarrow \mathbb{R}$ is the probability density function of the kinetic part in the momentum space \mathbb{R}^{3N} and $\pi_q : \Omega \rightarrow \mathbb{R}$ is the probability density function of the potential part in the configurational space Ω .

As we are interested in examining conformations/metastable sets, which are objects in configurational space, we will restrict ourselves to Ω :

“A conformation $C \subset \Omega$ will be identified with the particular metastable sub-ensemble $\mu_{C \times \mathbb{R}^{3N}}$ corresponding to the particular subset $C \times \mathbb{R}^{3N} \subset \Gamma$. Hence, for every position $q \in C$, the conformation contains all states with $q \in \Omega$ and arbitrary $p \in \mathbb{R}^{3N}$.”

In this sense, conformations/metastable sets contain no information on momenta and are determined in configurational space only. We are considering a reduced model in position space with a *reduced density* $\pi_q = \int_{\mathbb{R}^{3N}} \pi(q, p) \, dp$.

spatial/position
space
prob. dens.
fct.?

?

Macro States via Membership Functions

As the phase space and even the configurational space are very large, we aim to reveal the underlying discrete Markov State Model by group/cluster a collection of the micro states having the same or similar values in one observable. Such a collection of micro states will be called a *macro state*. For instance, that could be the states/observables “bound” or “unbound” for a receptor-ligand system.

We apply the function-based clustering method presented in section 2.3. We define macro states as overlapping partial densities, which can be identified as membership functions χ_1, \dots, χ_n . The membership functions $\chi_1, \dots, \chi_n : \Omega \rightarrow [0, 1]$ form a partition of unity, i.e.

$$\sum_{i=1}^n \chi_i(q) = 1. \quad (3.2)$$

By grouping micro states, the (corresponding) macro states yield *statistical weights*

$$w_i = \langle \chi_i, \mathbb{1} \rangle_\pi = \int_{\Omega} \chi_i(q) \pi_q(q) \, dq.$$

The statistical weight w_i corresponds to the “probability to be in conformation χ_i ”.

Transfer Operator

Each micro state $(q, p) \in \Gamma$ determines a *probability density function* $\Psi^{-\tau}(\cdot \mid (q, p))$ describing the possible evolutions of the system in configurational space Ω in time τ , being started at the initial state (q, p) . Weber[38] defines a transfer operator $\mathcal{P}(\tau) : L_{\pi_q}^{1,2}(\Omega) \rightarrow L_{\pi_q}^{1,2}(\Omega)$ for the propagation of (membership) functions via

$$\mathcal{P}(\tau)f(q) = \int_{\mathbb{R}^{3N}} \left(\int_{\Omega} f(\tilde{q}) \Psi^{-\tau}(\tilde{q} \mid (q, p)) \, d\tilde{q} \right) \pi_p(p) \, dp. \quad (3.3)$$

In this definition, the density function $\Psi^{-\tau}(\cdot \mid (q, p))$ can be interpreted as a transition function as defined in section 1.1. We have to notice that this transfer operator corresponds to the *backward operator* from section 1.2.

It is a *generalized* transfer operator in the sense that it includes deterministic as well as stochastic dynamical models. In order to describe deterministic dynamics, the density function $\Psi^{-\tau}$ has to be chosen as a Dirac delta function, since an initial state $(q(0), p(0))$ determines exactly the future states in configurational space.

It is important to remark that the transfer operator $\mathcal{P}(\tau)$ also defines a projected *Markov operator* $\bar{\mathcal{P}}(\tau)$ acting in configurational space Ω , see Weber[38], by

entropic inf.?
overlap =
good?

what is that
good for?
 e vs. $\mathbb{1}$

why not den-
sities?

?

propagator
sec 1.2
def

$$\overline{\mathcal{P}}(\tau) = \pi_q \circ \mathcal{P}(\tau) \circ (\pi_q)^{-1}, \quad (3.4)$$

which propagates density functions. The previous equation shows that the space of membership functions is connected to the space of density functions by multiplication with π_q . We will keep that relation in mind, but just use \mathcal{P} in the following.

As $\mathcal{P}(\tau)$ in (3.3) propagates **membership functions**, stationarity is characterized by the equation $e = \mathcal{P}(\tau)e$ for the constant function $e = 1$ in Ω . For the Markov operator $\overline{\mathcal{P}}(\tau)$ in (3.4) propagating **densities**, stationarity can be characterized by $\pi = \overline{\mathcal{P}}(\tau)\pi$, where π is the Boltzmann density. These two operators are *adjoint* operators. This can also be seen by the fact that a discretization of $\mathcal{P}(\tau)$ results in a matrix P_c , while a discretization of $\overline{\mathcal{P}}(\tau)$ will result in the transposed matrix P_c^T .

Boltzmann
dens. =
dist.?

Maybe: Properties of transfer operator for reversible Processes

Detailed Balance

$$\pi_q(\tilde{q}) \cdot \int_{\mathbb{R}^{3N}} \Psi^{-r}(q \mid (\tilde{q}, p)) \pi_p(p) \, dp = \pi_q(q) \cdot \int_{\mathbb{R}^{3N}} \Psi^{-r}(\tilde{q} \mid (q, p)) \pi_p(p) \, dp \quad (3.5)$$

Markov State Model for reversible Processes

For now, we consider a reversible process. Then due to the detailed balance condition (3.5), the corresponding transfer operator \mathcal{P} is **self-adjoint** and thus has a real spectrum, see theorem ?? (follows from linearity and self-adjointness) and $\sigma(\mathcal{P}) \subset [-1, 1]$ (since $\|\mathcal{P}f\|_{\pi_q} \leq \|f\|_{\pi_q}$). In order to apply the spectral approach from section 2.2, we assume that the **discrete spectrum** of the transfer operator \mathcal{P} has n **dominant eigenvalues** $1 = \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ which are all close to 1 and bounded away from the essential spectrum. The corresponding dominant eigenfunctions are denoted by $\mathcal{X} = \{\mathcal{X}_1, \dots, \mathcal{X}_n\}$ and therefore the eigenvalue problem is $\mathcal{P}(\tau)\mathcal{X} = \mathcal{X}\Lambda$, with the eigenvalue matrix $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$.

self-adj. wrt
 π_q

As we have seen in chapter 2, the number of metastable sets of a process can be determined by the number of dominant eigenvalues; i.e. we are going to create a Markov State Model on n states. The state space of this model should consist of the macro states of our Molecular System and its transition behaviour should be described via a $n \times n$ -transition matrix $P(\tau)$. In order to get from our continuous operator $\mathcal{P}(\tau)$ to a discrete matrix $P(\tau)$, we need at first to determine the size and shape of the membership functions χ_i . As described in section 2.3, this can be done by computing a

?

linear combination of the dominant eigenfunctions via

$$\chi_j(q) = \sum_{i=1}^N A_{ij} \mathcal{X}_i(q), \quad j = 1, \dots, n, \quad (3.6)$$

where $A = \{A_{ij}\}_{i,j=1,\dots,n}$ is the solution of PCCA+ (convex maximization problem). This choice of membership functions preserves Markovianity of the process when projecting. As a linear combination of eigenfunctions, the membership functions χ_i might have an overlap; they are not orthogonal!

PCCA+ only for finite/discrete state spaces?

Ref?

Galerkin Projection

Having computed the membership functions χ_i , we can project $\mathcal{P}(\tau)$ to a low-dimensional Markov State Model $P_c(\tau)$ by the Galerkin discretization

$$P_c(\tau) = G(\mathcal{P}(\tau)) = (\langle \chi, \chi \rangle_\pi)^{-1} (\langle \chi, \mathcal{P}(\tau) \chi \rangle_\pi). \quad (3.7)$$

In order to know about the quality of this model, we are interested in the iteration error under this projection. As mentioned in section 1.4, this error is zero if the Galerkin discretization of $(\mathcal{P}(\tau))^k$ is equal to the iteration $(P_c(\tau))^k$. In that case, the diagram 1.2 commutes. The following theorem shows that there is no discretization error under the projection (3.7), i.e. we have $(\mathcal{P}(\tau))^k = (P_c(\tau))^k$, which implies that Markovianity is preserved.

?

Theorem 3.1. (Weber [38, Theorem 2])

Let $\mathcal{P}(\tau)$ be the π_q -self-adjoint transfer operator defined in (3.3) with a set $\mathcal{X} = \{\mathcal{X}_1, \dots, \mathcal{X}_n\}$ of normalized eigenfunctions s.t. $\mathcal{P}(\tau)\mathcal{X} = \mathcal{X}\Lambda$ and a set of functions $\chi = \mathcal{X}A$ that is a linear combination of the eigenfunctions \mathcal{X} with a regular $n \times n$ -transformation matrix A from (3.6). Then the iteration error for the Galerkin discretization $P_c(\tau) = G(\mathcal{P}(\tau))$ in (3.7) vanishes.

Proof.

□

It follows that the above projection represents the correct dynamical long-time behaviour of the original process and that the matrix $P_c(\tau)$ is the correct Markov State Model. We can use the matrix representation $P_c(\tau) = S^{-1}T$ from theorem 1.17. Then S and T are stochastic matrices with

?

$$\begin{aligned} T &= D^{-1} \langle \chi, \mathcal{P}(\tau) \chi \rangle_\pi = D^{-1} A^T \Lambda A \quad \text{and} \\ S &= D^{-1} \langle \chi, \chi \rangle_\pi = D^{-1} A^T A, \end{aligned} \quad (3.8)$$

where $D = \text{diag}(w_1, \dots, w_n)$ is the diagonal matrix of statistical weights in (3.2).

Measuring the Rebinding Effect

Interpretation

Even though $P_c(\tau) := P(\tau)$ is the correct Markov State Model, it cannot be interpreted as a transition matrix, since the inverse matrix of S is not necessarily stochastic. The matrix $T = D^{-1}\langle\chi, \mathcal{P}\chi\rangle$ however can be interpreted as a transition matrix. Then the difference between $P_c(\tau)$ and T is given by

$$SP_c(\tau) = T.$$

Thus, the “disturbance” of ... can be measured by the matrix S . The more the matrix S differs from the identity matrix, the more the correct projection $P(\tau)$ differs from the transition matrix T . Thus, the rebinding effect can be measured by the matrix S . The trace of S is at most n . Optimizing $\text{trace}(S)$ is equivalent to optimizing the *crispness* of the conformations χ (Röblitz).

Infinitesimal Generator to transition rate matrix

Often (...) it is more convenient to consider/examine/investigate transition rate matrices instead of transition matrices/ infinitesimal generators instead of transfer operators. We can define the same/similar/analogous Galerkin Projection on the corresponding infinitesimal generator.

Conceptually, \mathcal{Q} is connected to the computation of transition rates.

The transfer operator $\mathcal{P}(\tau)$ defines a time-independent operator \mathcal{Q} via

$$\mathcal{Q} = \lim_{\tau \rightarrow 0} \frac{\mathcal{P}(\tau) - \mathcal{I}}{\tau},$$

which is the infinitesimal generator of \mathcal{P} :

Chapman

$$\mathcal{P}(\tau) = \exp(\tau\mathcal{Q}).$$

Weber[38] shows that such an infinitesimal generator exists for a discretization in terms of membership functions.

Since the eigenfunctions of \mathcal{Q} and \mathcal{P} are the same and their eigenvalues are related via $\exp(\xi_i) = \lambda_i$, we can apply the same Galerkin Projection for the infinitesimal generator as for the transfer operator in (3.7). We get a $n \times n$ -rate matrix

$$Q_c = A^{-1}\Xi A = (\langle\chi, \chi\rangle_\pi)^{-1}(\langle\chi, \mathcal{Q}\chi\rangle_\pi), \quad (3.9)$$

where Ξ is the diagonal matrix consisting of the n leading eigenvalues $0 = \xi_1 > \xi_2 \geq \dots \geq \xi_n$ of \mathcal{Q} and A is the transformation matrix of (3.6), which analogously transforms the eigenfunctions of \mathcal{Q} into membership functions of the macro states.

The matrix Q_c can be interpreted as a transition rate matrix.

?

3.3 Minimizing the Rebinding Effect

So far, we know that the matrix S from (3.8) gives a measure for the quantity of the rebinding effect, i.e. being close to the identity matrix implies a low rebinding, while high outer diagonal elements of S (overlap?) imply a high rebinding effect. But we don't know yet the meaning of this effect for the whole process. So we will at first set the rebinding effect (respectively the matrix S) in relation to the stability of a system/process. Afterwards, we will formulate an optimization problem in order to get a lower bound for the rebinding effect.

As the computation of eigenfunctions of a continuous operator \mathcal{Q} is an extensive task, we will assume in the further course that the transition rates can be measured experimentally. Thus, we will examine a given transition rate matrix Q_c .

Stability of the system/process in terms of determinants

If the eigenvalues ξ_i of Q_c are close to 0, then the macro states are very stable in the sense that the probability to stay inside of such a state is close to 1. The trace of Q_c corresponds to the sum of the dominant eigenvalues of \mathcal{Q} . Thus, we can measure the *stability* of the molecular system by considering $F := -\text{trace}(Q_c)$. If F is high, then the process is fast and less stable. If F is close to 0, then the process is slow and very stable. We want to set the indicator for stability F in relation to the matrices S and T from theorem 1.17 (matrix representation of Galerkin projection).

trace indep.
of A ?

Lemma 3.2. *If Q_c is a projected infinitesimal generator of a process and $P_c(\tau)$ the corresponding projected transfer operator with a matrix representation $P_c(\tau) = S^{-1}T$ from theorem 1.17, then the quantity $F := -\text{trace}(Q_c)$ can be measured by*

$$F = \tau^{-1}(\log(\det(S)) - \log(\det(T))).$$

Proof. We use the relation $\exp(\text{trace}()) = \det(\exp())$, the fact that Q_c “generates” $P_c(\tau)$, theorem 1.17 and multiplicativity of determinants to see that

$$\begin{aligned} F &= -\text{trace}(Q_c) \\ &= -\tau^{-1} \log(\exp(\text{trace}(\tau Q_c))) \\ &= -\tau^{-1} \log(\det(\exp(\tau Q_c))) \\ &= -\tau^{-1} \log(\det(P_c(\tau))) \\ &= \tau^{-1}(\log(\det(S)) - \log(\det(T))). \end{aligned}$$

see ..
why
 $\exp(\tau Q_c) = P_c(\tau)$?

□

Thus, both determinants of the stochastic matrices S and T influence the stability of the system, but in converse directions.

If $\det(T)$ is close to 1, then F is low and thus the process is stable/slow. If $\det(T)$ is close to 0, then the process is unstable/fast, since F is high. That makes sense because a high determinant of T can be interpreted as “good” metastability and thus corresponds to a slower/stable process, while a low determinant of T (“bad” metastability) makes the process faster/unstable.

assuming \det between 0 and 1?

On the other hand, if $\det(S)$ is close to 1 (almost unit matrix), then F is high, i.e. the process is unstable/fast. If $\det(S)$ is close to 0 (much overlap), then F is low, i.e. the process is stable/slow. That means that a higher overlap in S leads to a slower process! This relation is not as obvious at first sight.

As we figured out in section 3.2, the rebinding effect can be measured by the matrix S . If that matrix is close to the identity matrix, i.e. having a determinant close to 1, yields a low rebinding. Large deviation from the identity matrix, i.e. having higher outer diagonal elements, i.e. having a smaller determinant, results in high rebinding.

Combining these two properties of S , we can deduce that a high rebinding effect corresponds to a high stability of the system. This can be achieved by high outer diagonal elements/deviation from identity, i.e. by a high overlap of the membership functions(?).

Finding a lower bound for the rebinding effect

What is the meaning of the previous results/relation of S to stability of the system? On the one hand, the rebinding effect increases bindings(?). On the other hand, it increases the stability of a system. Thus, we are interested in a high/increased rebinding effect. We are going to compute a lower bound for it, in order to know how much rebinding there will be *at least*.

bad

In order to do so, let us first of all remember how S is determined. We were given a transfer operator \mathcal{P} which was projected onto a finite-dimensional state space via membership functions χ_i . These membership functions have been computed as a linear combination of the eigenfunctions with a regular matrix A . Thus, the choice of the matrix A determines S respectively the size of its determinant.

So far, A was assumed to be computed via PCCA+, i.e. such that the result is an optimal metastable decomposition/clustering. Now, we want to take into consideration the set of all possible, *feasible*, matrices A , to see if different choices result in a better/higher or worse/less rebinding.

Thus, we are going to formulate an optimization problem to find out which choice of A results in the lowest rebinding effect, measured by an *optimal matrix* S_{opt} , in order to know how much rebinding we are *guaranteed at least* for a given process. This problem is equivalent to finding the largest possible determinant of S .

“worst”

Optimization Problem (Maximizing determinant of S)

Since Q_c has the same eigenvalues as Q , the eigenvalue problem of Q_c is given by

$$Q_c X = X \Xi,$$

where the first column of X corresponds to the first eigenvector $X_1 := (1, \dots, 1)^T$. By (3.9), we see that A^{-1} is an eigenvector matrix of Q_c as well. Therefore, the columns of the matrix A^{-1} consists of multiples of the eigenvectors X_i . So we have

$$A^{-1} = \begin{pmatrix} 1 & & & \\ \vdots & \alpha_2 X_2 & \cdots & \alpha_3 X_3 \\ 1 & & & \end{pmatrix}$$

with $\alpha_2, \dots, \alpha_n \in \mathbb{R}$. We know from theorem 3.2 that a $\det(S)$ close to 1 results in a low rebinding effect. Thus, in order to find a lower bound for the rebinding effect, we try to maximize $\det(S)$, or equivalently minimize $|\det(S) - 1|$, since S is a stochastic matrix having 1 as largest possible determinant. The *objective function* of our optimization problem is then given by

$$\min_{\alpha_1, \dots, \alpha_n \in \mathbb{R}} |\det(S) - 1|, \quad (3.10)$$

where we have to include several *side constraints*. As the inverse matrix A^{-1} consists of linear combinations of eigenvectors X_i , we have to consider

$$\alpha_1 = 1 \quad \text{and} \quad A_{ij}^{-1} = \alpha_i X_{ij} \quad \forall i, j.$$

Furthermore, S is a stochastic matrix, see theorem 1.17, and its structure is given in terms of the linear transformation matrix A , so we have two further constraints

$$S = D^{-1} A^T A \quad \text{and} \quad S_{ij} \geq 0 \quad \forall i, j.$$

row-sum 1 included in formula?

A *feasible solution* of this optimization problem is a matrix S fulfilling all side constraints, but not necessarily being an optimum.

Interpretation

In the last section, we mentioned how the matrix S describes the *overlap* of the membership functions. For this reason, any feasible solution of the optimization problem (3.10) will be called a *real overlap matrix* S_{real} , while an actual optimum will be called an *optimal overlap matrix* S_{opt} . Clearly, we get $\det(S_{\text{real}}) \leq \det(S_{\text{opt}}) \leq 1$.

explain overlap ???

The real occurring rebinding effect is high if the determinant of S_{real} is low. Thus, a small determinant of S_{opt} increases the rebinding effect, while a large determinant

of S_{opt} gives us only few information about the quantity(?) of the rebinding effect, it could be large or small.

Unfortunately, for a reversible Q_c , the solution of optimization problem (3.10) gives us no information, as the following theorem shows.

Theorem 3.3. (Weber and Fackeldey[40, Theorem 1])

Let $Q_c \in \mathbb{R}^{n \times n}$ be a reversible matrix that stems from a clustering with positive definite overlap matrix S . Then there exists a matrix $A \in \mathbb{R}^{n \times n}$ in optimization problem (3.10) such that $\det(S_{\text{opt}}) = 1$.

Proof. It is enough to show that for a given reversible Q_c , we can find a matrix A fulfilling all constraints such that $S = D^{-1}A^T A$ is equal to the identity matrix I .

Assume that there is a regular matrix B , such that $Q_c = B^{-1}\Xi B$.

Since Q_c is reversible, we have $DQ_c = Q_c^T D$, see section 1.3, and thus

$$DB^{-1}\Xi B = B^T \Xi^T B^{-T} D.$$

Now let $C := B^{-T} D$, then we get

$$Q_c = \dots = C^{-1}\Xi C.$$

...

Have a real positive matrix $M = \text{diag}(m_1, \dots, m_n)$ and therefore a real positive diagonal matrix $\widetilde{M} = \text{diag}(\sqrt{m_1}, \dots, \sqrt{m_n})$. ???

...

Let $A := \widetilde{M}^{-1} B$. Show: A fulfills the constraints of (3.10), in order that S is a feasible matrix. Then

$$\begin{aligned} S &= D^{-1} A^T A = D^{-1} B^T \widetilde{M}^{-1} \widetilde{M}^{-1} B \\ &= D^{-1} B^T M^{-1} B \\ &= C^{-1} M^{-1} B \\ &= B^{-1} M M^{-1} B = I. \end{aligned}$$

Since all constraints of (3.10) are fulfilled, S is a feasible matrix with $\det(S) = 1$. \square

This theorem does **not** mean that a reversible process has no rebinding effect. It just means that for **every** reversible process, it is possible to find a clustering with no rebinding.

For instance, if we computed a clustering of a reversible process via PCCA+, then it could be the case that we have a low determinant of S , i.e. a high rebinding. But the optimization problem (3.10) gave us the lower bound of no rebinding, so it gave us no information about the rebinding of a particular/concrete clustering.

Linear Optimization Problem (Maximizing trace of S)

Now we present a different formulation of the above optimization problem (3.10). We will slightly change the objective function and turn the problem into a *linear optimization problem*. As a special case of the class of convex optimization problems, they have the nice property that any local optimum is also a global optimum.

and easier to solve?

We want to *minimize* the rebinding effect, i.e. we want to give a bound for how large the rebinding effect is *at least* (lower bound for rebinding effect). The closer the matrix S is to the identity matrix, the smaller is the rebinding effect. A matrix is close to the identity matrix, if its determinant is close to 1 or (equivalently) if its trace is close to n . So in order to compute the minimal rebinding effect, we can either maximize the determinant (get it as close as possible to 1) or maximize the trace of S (get it as close as possible to n), as the following theorem shows.

Theorem 3.4. *In optimization problem (3.10), we have $\det(S) \leq 1$ and $\text{trace}(S) \leq n$ with equality if and only if S is the identity/unit matrix.*

Proof.

□

So instead of maximizing $\det(S)$, we can also maximize $\text{trace}(S)$. In order to do so, let us first make some further observations about the eigenvectors of Q_c . We already found out that A^{-1} is a right eigenvector matrix of Q_c , with vectors being linear combinations of the eigenvectors X_i . Similarly, the matrix A is a *left* eigenvector matrix of Q_c , with row vectors being linear combinations of the eigenvectors Y_i . That fact can be expressed as

$$A = \tilde{U}Y^T = \begin{pmatrix} \tilde{\alpha}_1 Y_1 \\ \vdots \\ \tilde{\alpha}_n Y_n \end{pmatrix},$$

where each Y_i is a left eigenvector of Q_c (row vector) and the $\tilde{\alpha}_i \in \mathbb{R}$ are again some optimization parameters. The first eigenvector Y_i corresponds to the leading eigenvalue $\xi_1 = 0$ and is thus the stationary density of the process. The first row of A consists of the statistical weights of the clusters and therefore we have again $\tilde{\alpha}_i = 1$. With these notations we can write the new objective function as

$$\begin{aligned} \text{trace}(S) &= \text{trace}(D^{-1}A^T A) \\ &= \text{trace}(D^{-1}Y\tilde{U}^2 Y^T) \\ &= \sum_{i=1}^n \sum_{k=1}^n \tilde{\alpha}_k^2 \frac{y_{ik}^2}{y_{k1}}. \end{aligned}$$

The side constraints remain the same, i.e. $S_{ij} = \dots \geq 0$. Let $\beta = (\beta_1, \dots, \beta_n)$ with why $i \neq j$?

$\beta_i = \tilde{\alpha}_i^2$. Then the linear optimization problem of maximizing $\text{trace}(S)$ is given by

$$\max_{\beta} \sum_{k=1}^n \beta_k \left(\sum_{i=1}^n \frac{y_{ik}^2}{y_{k1}} \right), \quad (3.11)$$

fullfilling the side constraints

$$\beta_i \geq 0, \beta_1 = 1$$

and

$$\sum_{k=1}^n \beta_k y_{ik} y_{jk} \geq 0.$$

At first sight, this second formulation of the optimization problem might seem a bit more complex/confusing since we introduced several new matrices and variables. But in fact, the only change is that we maximize now the trace instead of the determinant (trace is easier to compute as it is just a sum). This formulation is better since, we have a *linear* program, which makes it easier to solve. And we have fewer constraints than before, because we merged some of the constraints into the objective function.

Let $B = \text{diag}(\beta_1, \dots, \beta_n)$. Then a solution β of (3.11) gives an optimal matrix $S_{\text{opt}} = D^{-1}YBY^T$ resulting in the smallest possible rebinding effect.

Conclusion

A nontrivial rebinding effect can be estimated only if the kinetics Q_c of a system is nonreversible.

why?

3.4 Approach for nonreversible processes

With the tools from section 2.4 (Schur Decomposition and G-PCCA+) we give an approach how this problem can be solved for nonreversible processes (NESS processes) using Schur Decomposition to get rid of the possibly nonreal eigenvalues, see Djurdevac et al[10](2016).

Transfer Operator

Have the transfer operator \mathcal{P} from (3.3) given, but from theorem 1.12 we know that the transfer operator of a nonreversible process is not self-adjoint.

Schur Decomposition

Applying a Schur Decomposition, we can create real eigenvalues from the possibly nonreal eigenvalues of the transfer operator.

Galerkin Projection

Now that we have real eigenvalues, we can apply the Galerkin Projection as usual.

detecting
dominant
cycles of the
process?
G-PCCA+?

4 Illustrative Examples

We want to apply the results from chapter 3 on some easy examples.

4.1 Transition Network Graph

4.2 Artificial (bivalent) binding Process

One can distinguish between a monovalent binding process and a multivalent binding process (see ...), where multivalent processes are often considered as having a better binding affinity (see..).

For the monovalent case, the mathematical modeling of its kinetics is well understood.

....

Whenever the receptor molecules are spatially preorganized, the corresponding binding process is denoted as multivalent.

(especialle bivalent or polyvalent case often observed in nature) These systems are of significant interest for pharmaceutical and technical applications. If the ligands are linked to each other in an appropriate way to match the preorganized receptor molecules and, thus, are also presented multivalently, then extremely high binding affinities are often observed.

So we consider here a bivalent process, as the the easiest multivalent case.

conformation	statistical weight	holding probability	life time (ps)
1	0.0001	0.9833	4.76
2	0.0003	0.9666	2.34
3	0.0041	0.9713	
4	0.0056	0.9987	

Table 4.1: relation of statistical weights to holding probabilities of the conformations

Conclusion

Bibliography

- [1] V. I. BOGACHEV, *Measure theory*, vol. 1, Springer Science & Business Media, 2007.
- [2] G. R. BOWMAN, V. S. PANDE, AND F. NOÉ, *An introduction to Markov state models and their application to long timescale molecular simulation*, vol. 797, Springer Science & Business Media, 2013.
- [3] J. H. BRANDTS, *Matlab code for sorting real schur forms*, Numerical linear algebra with applications, 9 (2002), pp. 249–261.
- [4] N. BROWN, *In Silico Medicinal Chemistry: Computational Methods to Support Drug Design*, no. 8 in Theoretical and Computational Chemistry Series, Royal Society of Chemistry, 2015.
- [5] B. R. CARÉ AND H. A. SOULA, *Impact of receptor clustering on ligand binding*, BMC Systems Biology, 5 (2011), p. 48.
- [6] J. D. CHODERA AND F. NOÉ, *Markov state models of biomolecular conformational dynamics*, Current opinion in structural biology, 25 (2014), pp. 135–144.
- [7] L.-T. DA, F. K. SHEONG, D.-A. SILVA, AND X. HUANG, *Application of markov state models to simulate long timescale dynamics of biological macromolecules*, in Protein Conformational Dynamics, Springer, 2014, pp. 29–66.
- [8] P. DEUFLHARD, W. HUISINGA, A. FISCHER, AND C. SCHÜTTE, *Identification of almost invariant aggregates in reversible nearly uncoupled markov chains*, Linear Algebra and its Applications, 315 (2000), pp. 39–59.
- [9] P. DEUFLHARD AND M. WEBER, *Robust perron cluster analysis in conformation dynamics*, Linear algebra and its applications, 398 (2005), pp. 161–184.
- [10] N. DJURDJEVAC CONRAD, M. WEBER, AND C. SCHÜTTE, *Finding dominant structures of nonreversible markov processes*, Multiscale Modeling & Simulation, 14 (2016), pp. 1319–1340.

- [11] K. FACKELDEY AND M. WEBER, *GenPCCA – markov state models for non-equilibrium steady states*, WIAS Report, 29 (2017), pp. 70–80.
- [12] B. GOLDSTEIN AND M. DEMBO, *Approximating the effects of diffusion on reversible reactions at the cell surface: ligand-receptor kinetics.*, Biophysical Journal, 68 (1995), p. 1222.
- [13] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, The John Hopkins University Press, Baltimore and London, 1996.
- [14] W. HUISINGA, *Metastability of markovian systems*, Ph.D. thesis, Freie Universität Berlin, (2001).
- [15] S. L. KALPAZIDOU, *Cycle representations of Markov processes*, vol. 28, Springer Science & Business Media, 2007.
- [16] T. KATO, *Perturbation Theory for Linear Operators*, Classics in Mathematics. Springer, 1995.
- [17] F. P. KELLY, *Reversibility and Stochastic Networks*, Wiley, 1979.
- [18] O. KNILL, *Probability and stochastic processes with applications*, Havard Web-Based, (1994).
- [19] D. A. MCQUARRIE, *Statistical Mechanics*, University Science Books, California, 2000.
- [20] S. P. MEYN AND R. L. TWEEDIE, *Markov Chains and Stochastic Stability*, Communications and Control Engineering Series. Springer, 1993.
- [21] A. NIELSEN, *Computation schemes for transfer operators*, Ph.D. thesis, Freie Universität Berlin, (2015).
- [22] B. ØKSENDAL, *Stochastic differential equations*, Springer, 2003.
- [23] J.-H. PRINZ, H. WU, M. SARICH, B. KELLER, M. SENNE, M. HELD, J. D. CHODERA, C. SCHÜTTE, AND F. NOÉ, *Markov models of molecular kinetics: Generation and validation*, The Journal of chemical physics, 134 (2011), p. 174105.
- [24] S. RÖBLITZ, *Statistical error estimation and grid-free hierarchical refinement in conformation dynamics*, PhD thesis, Freie Universität Berlin, 2009.
- [25] S. RÖBLITZ AND M. WEBER, *Fuzzy spectral clustering by PCCA+: application to markov state models and data classification*, Advances in Data Analysis and Classification, 7 (2013), pp. 147–179.

- [26] M. SARICH, *Projected transfer operators*, PhD thesis, Freie Universität Berlin, 2011.
- [27] M. SARICH, R. BANISCH, C. HARTMANN, AND C. SCHÜTTE, *Markov state models for rare events in molecular dynamics*, Entropy, 16 (2013), pp. 258–286.
- [28] C. SCHÜTTE, A. FISCHER, W. HUISINGA, AND P. DEUFLHARD, *A direct approach to conformational dynamics based on hybrid monte carlo*, Journal of Computational Physics, 151 (1999), pp. 146–168.
- [29] C. SCHÜTTE, W. HUISINGA, AND P. DEUFLHARD, *Transfer operator approach to conformational dynamics in biomolecular systems*, in Ergodic theory, analysis, and efficient simulation of dynamical systems, Springer, 2001, pp. 191–223.
- [30] C. SCHÜTTE AND M. SARICH, *Metastability and Markov State Models in Molecular Dynamics: Modeling, Analysis, Algorithmic Approaches*, vol. 24 of Courant Lecture Notes, American Mathematical Soc., 2013.
- [31] D. E. SHAW, R. O. DROR, J. K. SALMON, J. GROSSMAN, K. M. MACKENZIE, J. A. BANK, C. YOUNG, M. M. DENEROFF, B. BATSON, K. J. BOWERS, ET AL., *Millisecond-scale molecular dynamics simulations on anton*, in High performance computing networking, storage and analysis, proceedings of the conference on, IEEE, 2009, pp. 1–11.
- [32] D. E. SHAW, J. GROSSMAN, J. A. BANK, B. BATSON, J. A. BUTTS, J. C. CHAO, M. M. DENEROFF, R. O. DROR, A. EVEN, C. H. FENTON, ET AL., *Anton 2: raising the bar for performance and programmability in a special-purpose molecular dynamics supercomputer*, in Proceedings of the international conference for high performance computing, networking, storage and analysis, IEEE Press, 2014, pp. 41–53.
- [33] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Computer Science and Scientific Computing. Academic Press Boston, 1990.
- [34] K. STRØMGAARD, P. KROGSGAARD-LARSEN, AND U. MADSEN, *Textbook of Drug Design and Discovery*, CRC Press, 2002.
- [35] J. TOLLENAERE, *The role of structure-based ligand design and molecular modelling in drug discovery*, Pharmacy World and Science, 18 (1996), pp. 56–62.
- [36] G. VAUQUELIN, *Rebinding: or why drugs may act longer in vivo than expected from their in vitro target residence time*, Expert Opinion on Drug Discovery, 5 (2010), pp. 927–941.

- [37] M. WEBER, *Meshless methods in conformation dynamics*, Ph.D. thesis, Freie Universität Berlin, (2006).
- [38] M. WEBER, *A subspace approach to molecular markov state models via a new infinitesimal generator*, Habilitation thesis, Freie Universität Berlin, (2011).
- [39] M. WEBER, A. BUJOTZEK, AND R. HAAG, *Quantifying the rebinding effect in multivalent chemical ligand-receptor systems*, The Journal of Chemical Physics, 137 (2012), 054111.
- [40] M. WEBER AND K. FACKELDEY, *Computing the minimal rebinding effect included in a given kinetics*, Multiscale Modeling & Simulation, 12 (2014), pp. 318–334.
- [41] D. WERNER, *Funktionalanalysis*, Springer, 2006.