

# Finite element approach to clustering of multidimensional time series\*

Illia Horenko\*\*1

<sup>1</sup> Institut für Mathematik, Freie Universität Berlin  
Arnimallee 6, 14195 Berlin, Germany

**Key words** time series analysis, inverse problems, regularization, finite element method

**Subject classification** AMS: [62-07,62H30,62H25,65M60,60J10]

We present a new approach to clustering of time series based on a minimization of the *averaged clustering functional*. The proposed functional describes the mean distance between observation data and its representation in terms of  $\mathbf{K}$  abstract models of a certain predefined class (not necessarily given by some probability distribution). For a fixed time series  $x(t)$  this functional depends on  $\mathbf{K}$  sets of model parameters  $\Theta = (\theta_1, \dots, \theta_{\mathbf{K}})$  and  $\mathbf{K}$  functions of cluster affiliations  $\Gamma = (\gamma_1(t), \dots, \gamma_{\mathbf{K}}(t))$  (characterizing the affiliation of any element  $x(t)$  of the analyzed time series to one of the  $\mathbf{K}$  clusters defined by the considered model parameters). We demonstrate that for a fixed set of model parameters  $\Theta$  the appropriate Tykhonov-type regularization of this functional with some regularization factor  $\epsilon^2$  results in a minimization problem similar to a variational problem usually associated with one-dimensional non-homogeneous partial differential equation. This analogy allows us to apply the finite element framework to the problem of time series analysis and to propose a numerical scheme for time series clustering. We investigate the conditions under which the proposed scheme allows a monotone improvement of the initial parameter guess wrt. the minimization of the discretized version of the regularized functional. We also discuss the interpretation of the regularization factor in the Markovian case and show its connection to metastability and exit times.

The computational performance of the resulting method is investigated numerically on multi-dimensional test data and is applied to the analysis of multidimensional historical stock market data.

This is a preliminary version. Do not circulate!

## Introduction

Many application areas are characterized by the need to find some low-dimensional mathematical models for complex systems that undergo transitions between different phases. Such phases can be different circulation regimes in meteorology and climatology [1, 2, 3, 4], market phases in computational finance [5, 6] and molecular conformations in biophysics [7, 8, 9]. Regimes of this kind can sometimes not be directly observable (or "hidden") in the many dimensions of the system's degrees of freedom and can exhibit persistent or *metastable* behavior. If knowledge about the system is present only in the form of observation or measurement data, the challenging problem of identifying those metastable states together with the construction of reduced low-dimensional models becomes a problem of time series analysis and pattern recognition in high dimensions. The choice of the appropriate data analysis strategies (implying a set of method-specific assumptions on the analyzed data) plays a crucial role in correct interpretation of the available time series.

We present an approach to the identification of  $\mathbf{K}$  hidden regimes for an abstract class of problems characterized by a set of model-specific parameters  $\Theta = (\theta_1, \dots, \theta_{\mathbf{K}})$  and some positive bounded *model distance functional*  $g(x_t, \theta_i)$  describing a quality of data representation in terms of model  $i$ . We demonstrate how the hidden regimes can be obtained in form of cluster affiliations  $\Gamma = (\gamma_1(t), \dots, \gamma_{\mathbf{K}}(t))$  (characterizing the affiliation of any element of the analyzed time series to one of the  $K$  clusters defined

---

\* Supported by the DFG research center MATHEON "Mathematics for key technologies" in Berlin.

\*\* E-mail: horenko@math.fu-berlin.de

by the considered model parameters). In the *same way* as with other clustering methods described in the literature (like, f. e., K-Means [10, 11], fuzzy C-Means [12, 11] and fuzzy clustering with regression models (FCRM) algorithms [13]), and *in contrast* to Bayesian approaches (e.g., Gaussian Mixture Models (GMMs) [14, 15], Hidden Markov models (HMMs) [5, 16, 14, 17, 18, 19] or neuronal networks [20]), a key property of the presented numerical framework is that it does not impose any a priori probabilistic assumptions on the data. Later it will be explained how the additional probabilistic assumptions done a posteriori can help to identify the optimal number of cluster states. The proposed numerical scheme is based on the *finite element method (FEM)*, a technique widely used and studied in context of partial differential equations (PDEs). Application of the FEM in the context of time series analysis can potentially help to transfer the advanced numerical techniques currently developed in the PDE setting and allow for the construction of new *adaptive methods* of data analysis.

The remainder of this paper is organized in the following way: Sec. 1 presents a construction of the *regularized clustering functional* and demonstrates some examples of typical *model distance functionals*. Subsequently, a FEM discretization of the problem is derived in Sec. 2, a numerical optimization algorithm is presented and its properties are investigated. In Sec. 3 we give an interpretation of the *regularization factor* in terms of *metastable homogeneous Markov-jump processes*. Finally, numerical examples in Sec. 5 illustrate the use of the presented framework.

## 1 The averaged clustering functional and its regularization

### 1.1 Model distance functional

Let  $x(t) : [0, T] \rightarrow \Psi \subset \mathbf{R}^n$  be the observed  $n$ -dimensional time series. We look for  $\mathbf{K}$  *models* characterized by  $\mathbf{K}$  distinct sets of a priori unknown *model parameters*

$$\theta_1, \dots, \theta_{\mathbf{K}} \in \Omega \subset \mathbf{R}^d, \quad (1)$$

(where  $d$  is the dimension of a model parameter space) for the description of the observed time series. Let

$$g(x_t, \theta_i) : \Psi \times \Omega \rightarrow [0, \infty), \quad (2)$$

be a functional describing the *distance* from the observation  $x_t = x(t)$  to the *model*  $i$ . For a given *model distance functional* (2), under *data clustering* we will understand the problem of finding for each  $t$  a vector  $\Gamma(t) = (\gamma_1(t), \dots, \gamma_{\mathbf{K}}(t))$  called the *affiliation vector* (or vector of the *cluster weights*) together with model parameters  $\Theta = (\theta_1, \dots, \theta_{\mathbf{K}})$  which minimize the functional

$$\sum_{i=1}^{\mathbf{K}} \gamma_i(t) g(x_t, \theta_i) \rightarrow \min_{\Gamma(t), \Theta}, \quad (3)$$

subject to the constraints on  $\Gamma(t)$ :

$$\sum_{i=1}^{\mathbf{K}} \gamma_i(t) = 1, \quad \forall t \in [0, T] \quad (4)$$

$$\gamma_i(t) \geq 0, \quad \forall t \in [0, T], \quad i = 1, \dots, \mathbf{K}. \quad (5)$$

In the following, we will give three examples of the *model distance functional* (2) for three classes of cluster models: (I) *geometrical clustering*, (II) *Gaussian clustering* and (III) *clustering based on the essential orthogonal functions (EOFs)*.

**Example (I): Geometrical Clustering** One of the most popular clustering methods in multivariate data-analysis is the so-called *K-means algorithm* [21]. It is based on the iterative minimization of the distance from the data points to a set of  $K$  *cluster centers* which are recalculated in each iteration step. The affiliation to a certain cluster  $i$  is defined by the proximity of the observation  $x_t \in \Psi$  to the cluster

center  $\theta_i \in \Psi$ . In this case the *model distance functional* (2) takes the form of the square of the simple Euclidean distance between the points in  $n$  dimensions:

$$g(x_t, \theta_i) = \|x_t - \theta_i\|^2. \quad (6)$$

**Example (II): Gaussian Clustering** Another frequently used analysis algorithm is based on the identification of Gaussian sets in the analyzed data [16, 14]. It is assumed that the data  $x$  belonging to the same cluster  $i$  is distributed according to the multivariate normal distribution

$$p_i(x) = \sqrt{\det(2\pi\Sigma_i^{-1})} \exp\left(-0.5(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i)\right) \quad (7)$$

for all  $x \in \Psi$  with  $\theta_i = (\mu_i, \Sigma_i)$ ,  $\mu_i$  being the expectation value, and  $\Sigma_i$  the covariance matrix of  $p_i$ . In this case the *model distance functional* (2) can be expressed as a normed negative log-likelihood of (7):

$$g(x_t, \theta_i) = \|x_t - \mu_i\|_{\Sigma_i^{-1}}^2, \quad (8)$$

where  $\|\cdot\|_{\Sigma_i^{-1}}$  denotes the norm induced by the covariance matrix of the Gaussian distribution  $i^2$ .

**Example (III): EOF clustering** In many cases the dimensionality of the data  $x_t$  can be reduced to few *essential degrees of freedom* without significant loss of information. One of the most popular *dimension reduction* approaches used in applications is the method of *essential orthogonal functions* (EOFs) also known under the name of *principal component analysis* (PCA) [22]. As was demonstrated recently, it is possible to construct clustering methods based on the decomposition of data sets according to differences in their *essential degrees of freedom* allowing to analyze data of a very high dimensionality [17, 18, 19]. If the cluster  $i$  is characterized by a linear  $m$ -dimensional manifold ( $m \ll n$ ) of *essential degrees of freedom*, the respective model parameter is defined by the corresponding orthogonal projector  $\theta_i = \mathcal{T}_i \in \mathbf{R}^{n \times m}$  and the *model distance functional* (2) is given by the Euclidean distance between the original data  $x$  and its orthogonal projection on the manifold:

$$g(x_t, \theta_i) = \|x_t - \mathcal{T}_i \mathcal{T}_i^T x_t\|^2. \quad (9)$$

## 1.2 The averaged clustering functional and its regularization

Instead of solving the minimization problem (3) for each available element  $x_t \in \Psi, t \in [0, T]$  from the observed time series *separately*, one can approach all of the functional optimizations *simultaneously* and minimize the *averaged clustering functional*  $\mathbf{L}$ :

$$\mathbf{L}(\Theta, \Gamma) = \int_0^T \sum_{i=1}^{\mathbf{K}} \gamma_i(t) g(x_t, \theta_i) dt \rightarrow \min_{\Gamma, \Theta}, \quad (10)$$

subject to the constraints (1) and (4-5). The expression in (10) is similar to one that is typically used in the context of *Gaussian mixture models* (GMM) [14, 15] but is more general, since neither the function  $g(\cdot, \cdot)$  nor  $\Gamma(\cdot)$  have to be connected to some probabilistic models of the data (which is the case for *Gaussian mixture models*).

However, direct numerical solution of the problem (10) is hampered by the three following facts: (i) the optimization problem is *infinitely-dimensional* (since  $\Gamma(t)$  belongs to some not yet specified function class), (ii) the problem is *ill-posed* since the number of unknowns can be higher than the number of known parameters, and (iii) because of the non-linearity of  $g$  the problem is in general *non-convex* and

<sup>2</sup> It is important to mention that, in the context of many time series analysis methods, the *model distance functional* is not the *only quantity needed to formulate the numerical method*. For example, in the context of Bayesian methods (like GMMs and HMMs [16, 14]), the a priori known functional form of the probability distribution function (7) is also needed in each step of the algorithm. Therefore, in the following, we will draw a distinction between methods where the probabilistic assumptions should be implied a priori (like GMM/HMM) and clustering methods, which are based only on the notion of some (Euclidean, in many cases) *model distance functional*  $g(x_t, \theta_i)$  (like geometrical K-Means clustering methods [12]).

the numerical solution gained with some sort of *local minimization algorithm* depends on the initial parameter values [23].

One of the possibilities to approach the problems (i)-(ii) simultaneously is first to incorporate some *additional information* about the observed process (e.g., in the form of *smoothness assumptions* in space of functions  $\Gamma(\cdot)$ ) and then apply a finite Galerkin-discretization of this infinite-dimensional Hilbert space. For example, we can assume the *weak differentiability* of functions  $\gamma_i$ , i. e.:

$$|\gamma_i|_{\mathcal{H}^1(0,T)} = \|\partial_t \gamma_i(\cdot)\|_{\mathcal{L}_2(0,T)} = \int_0^T (\partial_t \gamma_i(t))^2 dt \leq C_\epsilon^i < +\infty, \quad i = 1, \dots, \mathbf{K}. \quad (11)$$

For a given observation time series, the above constraint limits the total number of transitions between the clusters and, as it will be demonstrated later in Section 3, is connected to the *metastability* of the *hidden process*  $\Gamma(t)$ .

Another possibility to incorporate the *a priori information* from (11) into the optimization is to modify the functional (3) and to write it in the *regularized form*

$$\mathbf{L}^\epsilon(\Theta, \Gamma, \epsilon^2) = \mathbf{L}(\Theta, \Gamma) + \epsilon^2 \sum_{i=1}^{\mathbf{K}} \int_0^T (\partial_t \gamma_i(t))^2 dt \rightarrow \min_{\Gamma, \Theta}. \quad (12)$$

This form of penalized regularization was first introduced by A. Tikhonov for solving ill-posed linear least-squares problems [24] and was frequently used for linear and non-linear regression analysis in context of statistics [25] and multivariate spline interpolation [26]. In contrast to the aforementioned applications of Tikhonov-type regularization to interpolation problems (where the regularization controls the smoothness of some non-linear functional approximation of the given data), the presented form of the regularized *averaged clustering functional* (12) has a completely different mathematical structure due to the linearity of the functional (10) wrt.  $\Gamma(t)$ . As it will be shown later, this specific formulation of the optimization problem with appropriate constrains on  $\gamma_i(t)$  allows to control the *metastability* (or *persistence*) of the assignment  $\Gamma(t)$  of the analyzed data to  $\mathbf{K}$  distinct a priori unknown clusters.

In the following, we will demonstrate a numerical approach to the optimization of this *regularized clustering functional*(12) subject to the constraints (1) and (4-5).

## 2 Finite elements approach to minimization of the regularized clustering functional

### 2.1 FEM-discretization

Let  $\{0 = t_1, t_2, \dots, t_{N-1}, t_N = T\}$  be a finite subdivision of the time interval  $[0, T]$ . We define a set of continuous functions  $\{v_1(t), v_2(t), \dots, v_N(t)\}$  with the *local support* on  $[0, T]$ , i. e.,  $v_1(t) \neq 0$  for  $t \in (t_1, t_2)$  (and zero elsewhere),  $v_k(t) \neq 0$  for  $t \in (t_{k-1}, t_{k+1})$ ,  $k = 2, \dots, N-1$  (and zero elsewhere),  $v_N(t) \neq 0$  for  $t \in (t_{N-1}, t_N)$  (and zero elsewhere). These functions are called *finite element basis* and there are lot of possible sets of such functions known from the literature on partial differential equations (PDEs) [27], like e.g., piecewise linear *hat functions* shown in Fig. 1.

Assuming that  $\gamma_i \in \mathcal{H}^1(0, T)$  we can write

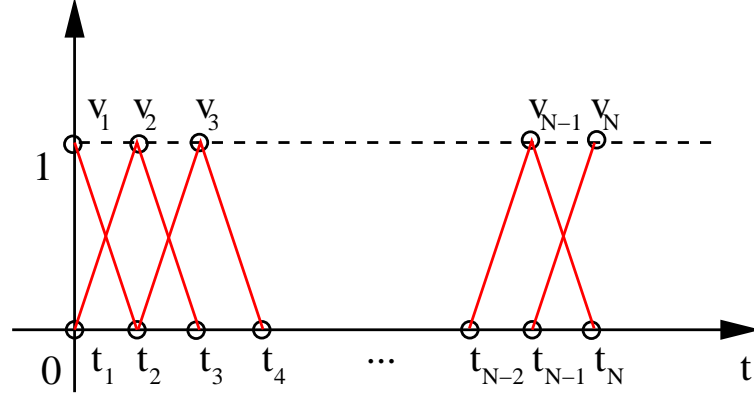
$$\begin{aligned} \gamma_i &= \tilde{\gamma}_i + \delta_N \\ &= \sum_{k=1}^N \tilde{\gamma}_i^{(k)} v_k + \delta_N, \end{aligned} \quad (13)$$

where  $\tilde{\gamma}_i^{(k)} = \int_0^T \gamma_i(t) v_k(t) dt$  and  $\delta_N$  is some *discretization error*. Substituting (13) in (12) we get

$$\mathbf{L}^\epsilon = \tilde{\mathbf{L}}^\epsilon + \mathcal{O}(\delta_N) \rightarrow \min_{\tilde{\gamma}_i, \Theta}, \quad (14)$$

where  $\tilde{\mathbf{L}}^\epsilon$  is a finite-dimensional version of the original functional (12):

$$\tilde{\mathbf{L}}^\epsilon = \sum_{i=1}^{\mathbf{K}} \int_0^T \left[ \tilde{\gamma}_i(t) g(x_t, \theta_i) + \epsilon^2 (\partial_t \tilde{\gamma}_i(t))^2 \right] dt. \quad (15)$$



**Fig. 1** Linear finite elements in one dimension.

After several obvious transformations and using of the locality of the finite element support we obtain:

$$\begin{aligned} \tilde{\mathbf{L}}^\epsilon = & \sum_{i=1}^{\mathbf{K}} \left[ \tilde{\gamma}_i^{(1)} \int_{t_1}^{t_2} v_1(t)g(x_t, \theta_i)dt + \sum_{k=2}^{N-1} \tilde{\gamma}_i^{(k)} \int_{t_{k-1}}^{t_{k+1}} v_k(t)g(x_t, \theta_i)dt + \right. \\ & + \tilde{\gamma}_i^{(N)} \int_{t_{N-1}}^{t_N} v_N(t)g(x_t, \theta_i)dt + \epsilon^2 \sum_{k=1}^{N-1} \left( \left( \tilde{\gamma}_i^{(k)} \right)^2 \int_{t_k}^{t_{k+1}} (\partial_t v_k(t))^2 dt - \right. \\ & \left. \left. - 2\tilde{\gamma}_i^{(k)}\tilde{\gamma}_i^{(k+1)} \int_{t_k}^{t_{k+1}} \partial_t v_k(t)\partial_t v_{k+1}(t)dt + \left( \tilde{\gamma}_i^{(k+1)} \right)^2 \int_{t_k}^{t_{k+1}} (\partial_t v_{k+1}(t))^2 dt \right) \right]. \quad (16) \end{aligned}$$

Denoting the vector of *discretized affiliations* to cluster  $i$  as  $\tilde{\gamma}_i = (\tilde{\gamma}_i^{(1)}, \dots, \tilde{\gamma}_i^{(N)})$ , vector of *discretized model distances* as

$$a(\theta_i) = \left( \int_{t_1}^{t_2} v_1(t)g(x_t, \theta_i)dt, \dots, \int_{t_{N-1}}^{t_N} v_N(t)g(x_t, \theta_i)dt \right), \quad (17)$$

and the symmetric tridiagonal *stiffness-matrix* of the finite element set as  $\mathbf{H}$

$$\mathbf{H} = \begin{pmatrix} \int_{t_1}^{t_2} v_1^2(t)dt & \int_{t_1}^{t_2} v_1(t)v_2(t)dt & 0 & \dots & 0 \\ \int_{t_1}^{t_2} v_1(t)v_2(t)dt & \int_{t_2}^{t_3} v_2^2(t)dt & \int_{t_2}^{t_3} v_2(t)v_3(t)dt & \dots & 0 \\ \dots & \dots & \dots & \ddots & 0 \\ 0 & 0 & 0 & \dots & \int_{t_{N-1}}^{t_N} v_N^2(t)dt \end{pmatrix}. \quad (18)$$

With the help of (17-18) we can re-write (16) as

$$\tilde{\mathbf{L}}^\epsilon = \sum_{i=1}^{\mathbf{K}} [a(\theta_i)^{\mathbf{T}}\tilde{\gamma}_i + \epsilon^2\tilde{\gamma}_i^{\mathbf{T}}\mathbf{H}\tilde{\gamma}_i] \rightarrow \min_{\tilde{\gamma}_i, \Theta}, \quad (19)$$

subject to (1), the discretized version of equality constraints (4)

$$\sum_{i=1}^{\mathbf{K}} \tilde{\gamma}_i^{(k)} = 1, \quad \forall k = 1, \dots, N, \quad (20)$$

and the inequality constraints (5)

$$\tilde{\gamma}_i^{(k)} \geq 0, \quad \forall k = 1, \dots, N, \quad i = 1, \dots, \mathbf{K}. \quad (21)$$

## 2.2 Numerical method and monotonicity conditions

The minimization problem (19-21), for a fixed set of *cluster model parameters*  $\Theta$  reduces to a quadratic optimization problem with linear constraints which can be solved by standard tools of quadratic programming (QP) like, e.g., the *ellipsoid methods* or the *interior point methods* that converge in polynomial time [28, 29, 23]. If, in addition, it is possible to minimize the problem (19-21) wrt. parameters  $\Theta$  for a *fixed* set of *discretized cluster affiliations*  $\bar{\gamma}_i$ , we can split the original optimization problem in two consecutive parts repeated in each iteration step of the following algorithm:

### Algorithm.

*Setting of optimization parameters and generation of initial values:*

- Set the number of clusters  $\mathbf{K}$ , regularization factor  $\epsilon^2$ , finite discretization of the time interval  $[0, T]$ , and the optimization tolerance TOL
- Set the iteration counter  $j = 1$
- Choose random initial  $\bar{\gamma}_i^{[1]}, i = 1, \dots, \mathbf{K}$  satisfying (20-21)
- Calculate  $\Theta^{[1]} = \arg \min_{\Theta} \tilde{\mathbf{L}}^\epsilon(\Theta, \bar{\gamma}_i^{[1]})$  subject to (1)

*Optimization loop:*

- do**
- Compute  $\bar{\gamma}^{[j+1]} = \arg \min_{\bar{\gamma}} \tilde{\mathbf{L}}^\epsilon(\Theta^{[j]}, \bar{\gamma})$  satisfying (20-21) applying QP
  - Calculate  $\Theta^{[j+1]} = \arg \min_{\Theta} \tilde{\mathbf{L}}^\epsilon(\Theta, \bar{\gamma}_i^{[j+1]})$
  - $j := j + 1$
- while**  $\left| \tilde{\mathbf{L}}^\epsilon(\Theta^{[j+1]}, \bar{\gamma}_i^{[j+1]}) - \tilde{\mathbf{L}}^\epsilon(\Theta^{[j]}, \bar{\gamma}_i^{[j]}) \right| \geq \text{TOL}.$

Note that the solution of the QP problem wrt.  $\bar{\gamma}$  will always *exactly* satisfy the constraints (20-21) guaranteeing that the resulting posterior probabilities (or cluster affiliations)  $\gamma_i(t)$  all add to one for each  $t$ . This will guarantee the mathematical coherence of the implied affiliations  $\gamma_i(t)$ , cf. [30].

A solution of the problem  $\Theta^{[j+1]} = \arg \min_{\Theta} \tilde{\mathbf{L}}^\epsilon(\Theta, \bar{\gamma}_i^{[j+1]})$  can be calculated by, e.g., re-writing the problem in Lagrangian form and applying the Newton–method. Moreover, in many cases this problem can even be solved analytically (this depends on the form of the *model distance functional*). For all 3 examples of  $g(x_t, \theta)$  presented above this can be done. Note that due to the non-linearity of  $g(\cdot, \theta)$  wrt.  $\theta$ , the minimized functional  $\tilde{\mathbf{L}}^\epsilon$  is non-convex. This feature will prohibit us to use the standard results of convex optimization theory to show the *convergence* of the presented algorithm. The following theorem describes the conditions under which the above algorithm will *monotonously minimize* the *energy* (19).

**Theorem 2.1** *Let for a given observed time series  $x(t) : [0, T] \rightarrow \Psi \subset \mathbf{R}^n$  the model distance functional  $g$  be chosen such that (2) is fulfilled,  $\Psi$  and  $\Omega$  are compact,  $g(x_t, \cdot)$  is continuously differentiable function of  $\theta$  and*

$$\frac{\partial}{\partial \Theta} \tilde{\mathbf{L}}^\epsilon(\Theta^*, \bar{\gamma}) = 0 \quad (22)$$

*has a solution  $\Theta^* = (\theta_1^*, \dots, \theta_{\mathbf{K}}^*), \theta_{i^*} \in \Omega$  for any fixed  $\bar{\gamma}$  satisfying (20-21) and  $\frac{\partial^2}{\partial \Theta^2} \tilde{\mathbf{L}}^\epsilon(\Theta^*, \bar{\gamma})$  exists and is positive definite. Then for any  $\epsilon^2 \geq 0$  and any finite non-negative finite elements set  $\{v_1(t), v_2(t), \dots, v_N(t)\} \in \mathcal{L}_2(0, T)$  such that the respective stiffness-matrix  $\mathcal{H}$  is positive semidefinite, the above algorithm is monotone, i. e., for any  $j \geq 1$*

$$\tilde{\mathbf{L}}^\epsilon(\Theta^{[j+1]}, \bar{\gamma}_i^{[j+1]}) \leq \tilde{\mathbf{L}}^\epsilon(\Theta^{[j]}, \bar{\gamma}_i^{[j]}). \quad (23)$$

**Proof.** Since  $\epsilon^2 \mathbf{H}$  is positive semidefinite and  $0 \leq g(x_t, \theta) \leq \bar{g} < +\infty$ , for any fixed  $\theta \in \Omega$  the functional  $\tilde{\mathbf{L}}^\epsilon$  is convex. Moreover, (2) implies that  $\tilde{\mathbf{L}}^\epsilon$  is bounded from below and constraints (20-21) define a non-empty closed convex domain. In this case the problem  $\bar{\gamma}^{[j+1]} = \arg \min_{\bar{\gamma}} \tilde{\mathbf{L}}^\epsilon(\Theta^{[j]}, \bar{\gamma})$

satisfying (20-21) has a *global minimizer*  $\bar{\gamma}_i^{[j+1]}$ , in particular

$$\tilde{\mathbf{L}}^\epsilon \left( \Theta^{[j]}, \bar{\gamma}_i^{[j+1]} \right) \leq \tilde{\mathbf{L}}^\epsilon \left( \Theta^{[j]}, \bar{\gamma}_i^{[j]} \right). \quad (24)$$

Moreover, due to (22) and since the Hesse-matrix  $\frac{\partial^2}{\partial \Theta^2} \tilde{\mathbf{L}}^\epsilon \left( \Theta^*, \bar{\gamma}_i^{[j+1]} \right)$  exists and is positive-definitive, the solution of  $\Theta^{[j+1]} = \arg \min_{\Theta} \tilde{\mathbf{L}}^\epsilon \left( \Theta, \bar{\gamma}_i^{[j+1]} \right)$  subject to (1) exists and

$$\tilde{\mathbf{L}}^\epsilon \left( \Theta^{[j+1]}, \bar{\gamma}_i^{[j+1]} \right) \leq \tilde{\mathbf{L}}^\epsilon \left( \Theta^{[j]}, \bar{\gamma}_i^{[j+1]} \right). \quad (25)$$

Finally, (24) together with (25) results in (23).  $\square$

### 3 The Markovian case: regularization factor and metastability

The proposed numerical method results in a *local improvement of the energy* (19). The minimization problem was obtained by a finite element discretization of the continuous *regularized clustering problem* (12) under conditions (4-5). However, it is not a priori clear what is the connection between the discrete solution we obtain with the above algorithm and the minimizer of the original *averaged clustering functional* (10). There are two main questions to be answered: (i) what is the *discretization error*  $\delta_N$  introduced on the way from (12) to (19) and how does it influence the quality of the resulting minimizers, and (ii) what is the influence of the regularization factor  $\epsilon^2 \sum_{i=1}^{\mathbf{K}} \int_0^T (\partial_t \gamma_i(t))^2 dt$ .

Concerning the first question, it seems clear that with increasing number  $N$  of time discretization points the error  $\delta_N$  from (13) will decrease and the overall difference between the continuous and discretized versions of the regularized functional will be getting smaller. For a rigorous mathematical justification of this feature and for the estimation of the discretization error, one can apply the theory developed for partial differential equations. This is a matter of future research. Here we would only like to mention the fact that in practical applications the observation time series are almost always available only in a discrete form (since the measurements of the real life processes can be typically acquired only at some discrete moments in time). This means that one actually starts with a discretized problem and therefore there is an upper limit for  $N$  given by the number of times the process was observed. However, we want to keep the continuous representation of the optimization problem in order to be able to construct the *adaptive FEM* scheme in future. Control over the *discretization error* in (13) will allow implementation of the adaptivity exactly in the same way as it is done in the theory of PDEs [27].

Concerning the influence of the regularization factor, it is intuitively clear that the penalization of the derivative norm in form of regularization (12) or in form of the constraint (11) has a *smoothing effect*. More specifically, in the case of a piecewise-constant function  $\Gamma$ , regularization will result in restriction of the number of transitions between the clusters. This means that the *cluster affiliation functions*  $\gamma_i$  obtained by the optimization of the *regularized functional* will be more and more *metastable* for increasing  $\epsilon^2$  or  $C_{\epsilon^2}$ , i. e., the observed process will stay more and more time in the respective identified cluster before making a transition to the next identified cluster. Metastability is associated with the so called *mean exit times*  $\tau_i^{\text{exit}}$  describing the mean time a process will stay in the state  $i$  before leaving to some other state. In the context of *discrete homogeneous Markovian processes*, the *mean exit time*  $\tau_i^{\text{exit}}$  can be quantified with the help of the *transition matrix*  $\mathbf{P}$  [31]:

$$\tau_i^{\text{exit}} = \frac{\Delta_t}{1 - \mathbf{P}_{ii}}, \quad (26)$$

where  $\Delta_t$  is the discrete time step of the Markov chain and  $\mathbf{P}_{ii}$  is the Markovian probability to stay in the same state  $i$  after one time step  $\Delta_t$ .

We will investigate the effect of regularization for a discretized optimization problem. As a finite element basis we take the *linear finite elements* shown in Fig. 1 on an equidistant time grid with step



$\Delta_t$ :

$$v_k(t) = \begin{cases} \frac{t-t_k}{\Delta_t} & 2 \leq k \leq N-1, t \in [t_{k-1}, t_k], \\ \frac{t_{k+1}-t}{\Delta_t} & 2 \leq k \leq N-1, t \in [t_k, t_{k+1}], \\ \frac{t_2-t}{\Delta_t} & k=1, t \in [t_1, t_2] \\ \frac{t-t_{N-1}}{\Delta_t} & k=N, t \in [t_{N-1}, t_N]. \end{cases} \quad (27)$$

The following theorem explains a connection between the *regularization factor* and *metastability* of the resulting clustering in the Markovian case.

**Theorem 3.1** *Let  $\Gamma^\epsilon(t) = (\gamma_1^\epsilon, \dots, \gamma_{\mathbf{K}}^\epsilon)$  be a solution of the optimization problem (19-21), (1) resulting from application of the positive linear finite elements discretization  $v_l(t)$  (27) with a constant time step  $\Delta_t$  and*

$$C_0^i(\epsilon) = \sum_{l=1}^N [\gamma_{i,l}^\epsilon] \quad (28)$$

$$C_1^i(\epsilon) = [\gamma_i^\epsilon]^T \mathbf{H} [\gamma_i^\epsilon], \quad i = 1, \dots, \mathbf{K}, \quad (29)$$

where  $[\cdot]$  is a component-wise roundoff operation towards the nearest integer. If the values  $\gamma_i(t_k) = \sum_{l=1}^N [\gamma_{i,l}^\epsilon] v_l(t_k)$ ,  $i = 1, \dots, \mathbf{K}$  of respective cluster affiliation function are considered as an output of a time-discrete homogeneous Markov-jump process with time step  $\Delta_t$  (where  $t_1, \dots, t_N$  is the subdivision of  $[0, T]$ ), then the respective mean exit times  $\tau_i^{exit}$  are

$$\tau_i^{exit} = \frac{C_0^i(\epsilon)}{C_1^i(\epsilon)} \Delta_t, \quad i = 1, \dots, \mathbf{K}. \quad (30)$$

*Proof.* Since  $v_k(t)$  has the form of (27), the *stiffness-matrix*  $\mathbf{H}$  in (19) will be symmetric and tridiagonal with  $2/\Delta_t$  on the main diagonal,  $-1/\Delta_t$  on both secondary diagonals and zero elsewhere. Therefore we can write:

$$C_1^i(\epsilon) = \frac{1}{\Delta_t} \sum_{l=1}^{\mathbf{K}} \left( [\gamma_{i,(l+1)}^\epsilon] - [\gamma_{i,l}^\epsilon] \right)^2. \quad (31)$$

Moreover

$$C_1^i(\epsilon) \leq C_0^i(\epsilon) \leq N, \quad (32)$$

since  $[\gamma_{i,l}^\epsilon]$  can only take values 0 and 1. Therefore,  $(C_0^i(\epsilon) - C_1^i(\epsilon))$  is the number of times the observation process stayed in cluster  $i$  without going elsewhere in the next step and  $C_0^i(\epsilon)$  is the total number of times the state  $i$  was visited. The *maximum log-likelihood* estimate of the respective homogeneous Markovian probability is

$$\mathbf{P}_{ii} = \frac{(C_0^i(\epsilon) - C_1^i(\epsilon))}{C_0^i(\epsilon)}, \quad i = 1, \dots, \mathbf{K}, \quad (33)$$

and the corresponding *maximum log-likelihood* estimate of the *mean exit time* is given by the expression (30).  $\square$

The above theorem demonstrates a connection between the effect of *regularization* and *metastability* of the resulting clustering when the finite element discretization with uniform time step  $\Delta_t$  is interpreted as the output of a homogeneous Markov-jump process on the same time scale. It is intuitively clear that with growing  $\epsilon^2$  the energy-norm  $\|\cdot\|_{\mathcal{H}^1(0,T)}$  of the resulting optimal vector  $\gamma_i^\epsilon$  will get smaller. It means that according to (30), the respective mean exit times are increasing and the corresponding cluster decomposition becomes more metastable if the  $\mathcal{H}^1(0, T)$  half-norm of the solution is changing faster than its  $\mathcal{L}^1(0, T)$  norm. As it will be demonstrated later, this property can be used to determine the optimal value of the *regularization parameter*  $\epsilon^2$  in the context of the standard L-curve approach to Tikhonov-regularized problems [32, 33].



## 4 A posteriori probabilistic assumptions: choice of $\mathbf{K}$ and $\epsilon^2$

**Estimation of optimal  $\mathbf{K}$  dependent on  $\epsilon^2$ :** The quality of the resulting clustering and reliability of the model parameters  $\Theta$  in any specific case will be very much dependent on the original data, especially on the length of the available time series. The shorter the observation sequence, the bigger is the uncertainty of the resulting parameters. The same is true if the number  $\mathbf{K}$  of clusters is being increased for fixed length of the observed time series: the bigger  $\mathbf{K}$ , the higher will be the uncertainty for each of the states. Therefore identified cluster states should be distinguishable in some sense. The *upper bound for the number of statistically distinguishable cluster states* for each value of  $\epsilon^2$  can be algorithmically estimated in the following way: starting with some a priori chosen (big)  $\mathbf{K}$  one solves the optimization problem (12) for a fixed value of  $\epsilon^2$  and calculates the confidence intervals of the resulting local parameters  $\theta_i, i = 1, \dots, \mathbf{K}$ . This can be done implying some probabilistic assumptions, f. e., *independence and identity of distributions* (i.i.d.) or Gaussianity of cluster distributions, for each of the identified clusters a posteriori. If two of the estimated sets of parameters  $\theta_i$  and  $\theta_j$  have the confidence intervals that are overlapping in all components, this means that respective clusters  $i$  and  $j$  are *statistically indistinguishable* and the whole procedure must be repeated for  $\mathbf{K} = \mathbf{K} - 1$ . If at certain point all of the clusters are *statistically distinguishable* the procedure is stopped and  $\mathbf{K}_{max}(\epsilon^2) = \mathbf{K}$ .

Alternatively, if the a posteriori analysis of the cluster affiliation functions  $\gamma_i(t)$  reveals the Markovian nature of the hidden transition process [34], spectral theory of Markov chains can be applied [35]. It can help to determine the number  $\mathbf{K}_{max}$  of *metastable* states from the number of the dominant eigenvalues in the so called *Perron cluster* [35]. For example, the *Perron cluster - cluster analysis (PCCA)* [36] can be used for the a posteriori analysis of  $\gamma_i(i)$  to find a lower bound for  $\mathbf{K}_{max}$ .

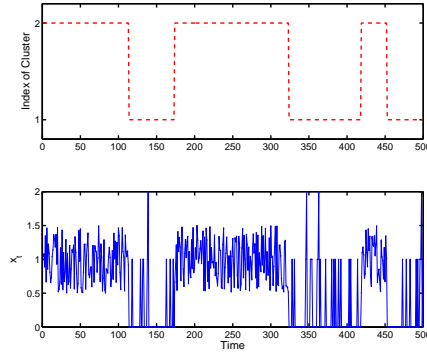
As was already mentioned above, in contrast to GMM/HMM approaches [14, 15] and similar to standard clustering techniques (like K-Means clustering [12]), the FEM-based clustering procedure itself does not rely on any a priori probabilistic assumptions. Both of the ways of determining the  $\mathbf{K}_{max}$  described above rely on the probabilistic assumptions about the observed or hidden data made a posteriori, i. e. after performing the FEM-based minimization of the functional (12) for certain fixed values of  $\mathbf{K}, \epsilon^2$ . Combination of both above approaches can help to find the *optimal* number  $\mathbf{K}$  of clusters in each specific application. This will be demonstrated on the numerical example later in the text.

**Estimation of optimal  $\epsilon^2$ :** As was already mentioned above, linearity of the original *averaged clustering functional* (10) wrt.  $\gamma_i(t)$  for *fixed* values of model distance parameters  $\theta_i$  can help to apply the standard instruments from the theory of ill-posed linear problems, like Morozov discrepancy principle [37, 38] and L-curve approach [32, 33] to identify the optimal value of the regularization parameter  $\epsilon^2$ . For example, in context of the L-curve method, optimal value of the Tikhonov-parameter is determined as the edge-point (or the point of maximal curvature) on the two-dimensional plot. This plot depicts a dependence of the residuum-norm of the solution from the norm of the regularized solution (calculated for different values of regularization parameter). To apply this idea in the context of problem (12), a value of the functional (10) and the  $\mathcal{H}_1(0, T)$  half-norm of the respective solution can be taken to investigate the L-curve behavior. As was shown in Theorem 3.1, there is also a connection between the  $\mathcal{H}_1(0, T)$  half-norm and *mean exit times*  $\tau_i$  (measuring persistence of clustering in Markovian case). As it will be demonstrated later in the numerical example, information about the exit times can also be used to determine the edge point of the L-curve.

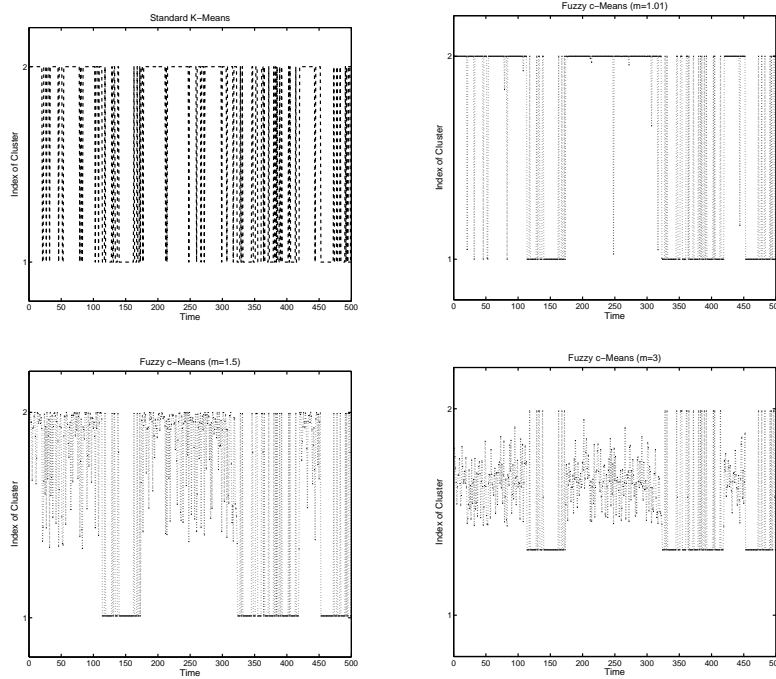
## 5 Numerical examples

To investigate the proposed framework numerically, three examples will be presented. (i) In the first numerical example, the application to the one-dimensional test model system is presented and comparison with standard geometric clustering methods (*K-means* and *fuzzy C-Means*) is performed. (ii) In the second numerical example, a multidimensional test case is considered where the time series is constructed from two multivariate non-Gaussian distributions (which prohibits application of standard GMM/HMM methods) with equal expectation values and different covariances (which prohibits the application of standard K-Means and fuzzy C-Means clustering algorithms). (iii) In the third example,

application of the presented FEM-clustering methodology to the analysis of 1106 NASDAQ stock prices in year 2007 is presented.



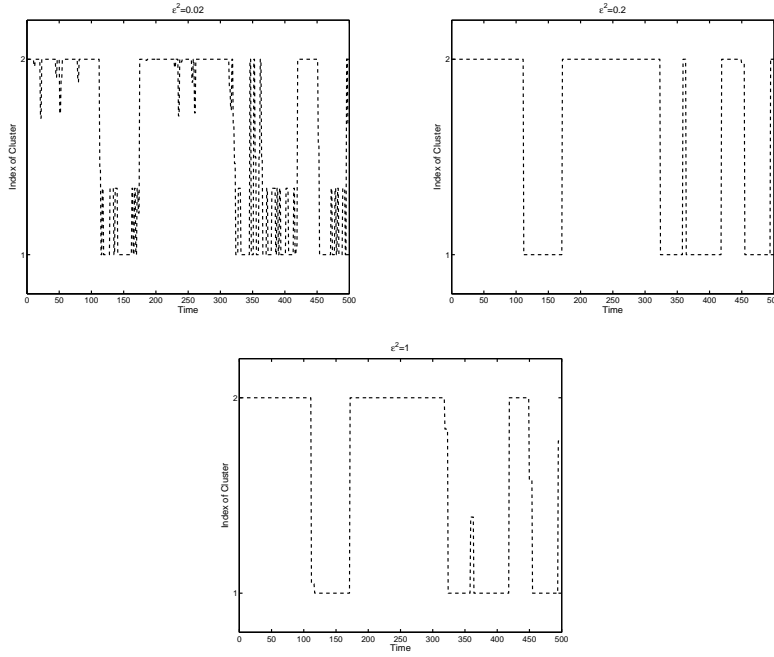
**Fig. 2** Hidden discrete process  $Y_t$  (upper panel) switching between two model distributions, Poisson distribution with expectation value 0.5 for the cluster 1 and equally distributed random variable in interval  $[0.5, 1.5]$  for the cluster 2. The resulting time series  $x_t$  is shown in the lower panel of the plot.



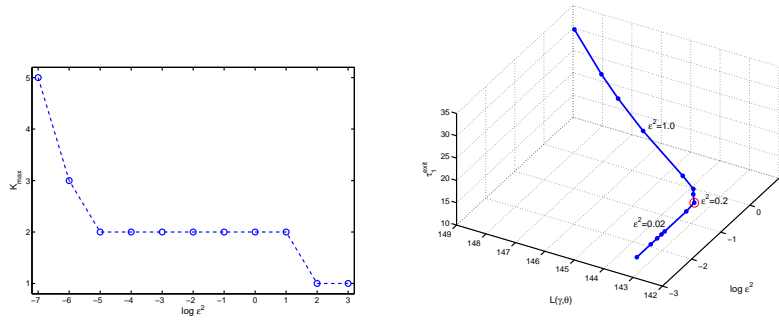
**Fig. 3** Cluster affiliations resulting from application of the standard clustering algorithms to the time series from the lower panel of Fig. 2: (i) K-Means clustering procedure (upper left panel) and (ii) fuzzy C-Means algorithm for different values of the fuzzyfier parameter  $m$ . The number of clusters is set to  $\mathbf{K} = 2$  in each case. The algorithms are initialized with 100 different randomly chosen initial cluster parameters to avoid trapping in the local optimum, results with the lowest value of the cluster distance functional are shown.

### 5.1 One-dimensional test data: choice of $\mathbf{K}$ , $\epsilon^2$ and comparison to standard methods

First we apply the method to a one dimensional test time series generated by the discrete hidden process  $Y_t$  given in the upper panel of Fig. 2. Each of the two hidden states is represented by a certain

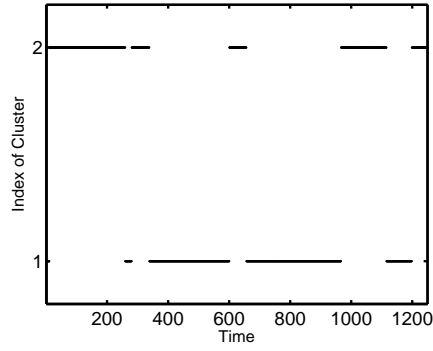


**Fig. 4** Cluster affiliations  $\gamma_1(t)$  resulting from application of the FEM-based minimization of the regularized averaged clustering functional (12) to the time series from the lower panel of Fig. 2 (with *model distance functional* of the form (6)). Three graphic panels correspond to affiliations obtained for three different values of the regularization parameter  $\epsilon^2 = 0.02, 0.2, 1.0$  ( $\mathbf{K} = 2, N = 500$ ). The algorithm is initialized with 100 different randomly chosen initial cluster parameters to avoid trapping in the local optimum, results with the lowest value of the functional (12) are shown.

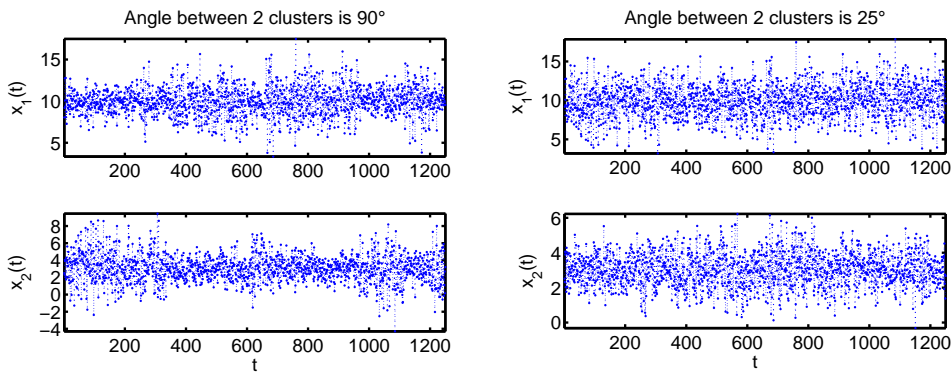


**Fig. 5** Left panel: dependence between the choice of a regularization parameter  $\epsilon^2$  and the maximal number of *statistically distinguishable cluster states*  $\mathbf{K}_{max}$  (for any fixed  $\epsilon^2$ ,  $\mathbf{K}_{max}$  is determined according to the procedure described in Sec. 4). Number of linear finite element functions  $N$  is chosen to be equal to the length of the analyzed time series ( $N = T = 500$ ). Right panel: dependence between  $\epsilon^2$ , mean exit time and the respective value of the original averaged model distance functional (10) in the optimum of the regularized problem (12).

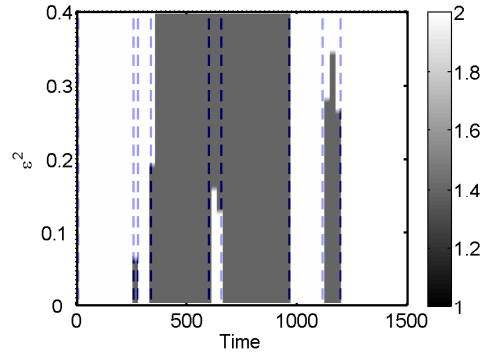
probability distribution of the observed process  $x_t$ :  $Y_t = 1$  corresponds to a Poisson distribution of  $x_t$  with expectation value 0.5 and  $Y_t = 2$  is associated with equally distributed random variables  $x_t$  in interval  $[0.5, 1.5]$ . Resulting time series with 500 elements is shown in the lower panel of Fig. 2. As can be seen from the Fig. 2, resulting time series  $x_t$  exhibits strongly overlapping non-Gaussian processes of very different probabilistic nature (one being Poisson- and one being equally-distributed).



**Fig. 6** Hidden discrete process switching between data-clusters.



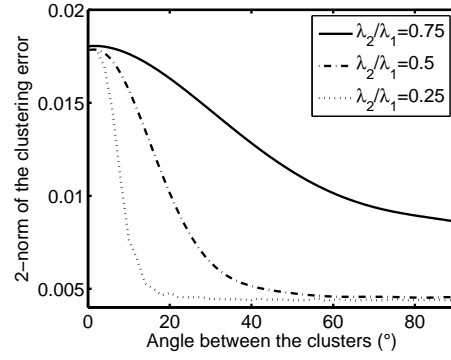
**Fig. 7** Time series in "Gaussian" degrees of freedom  $x_1, x_2$  generated by the hidden switching process from Fig. 6.



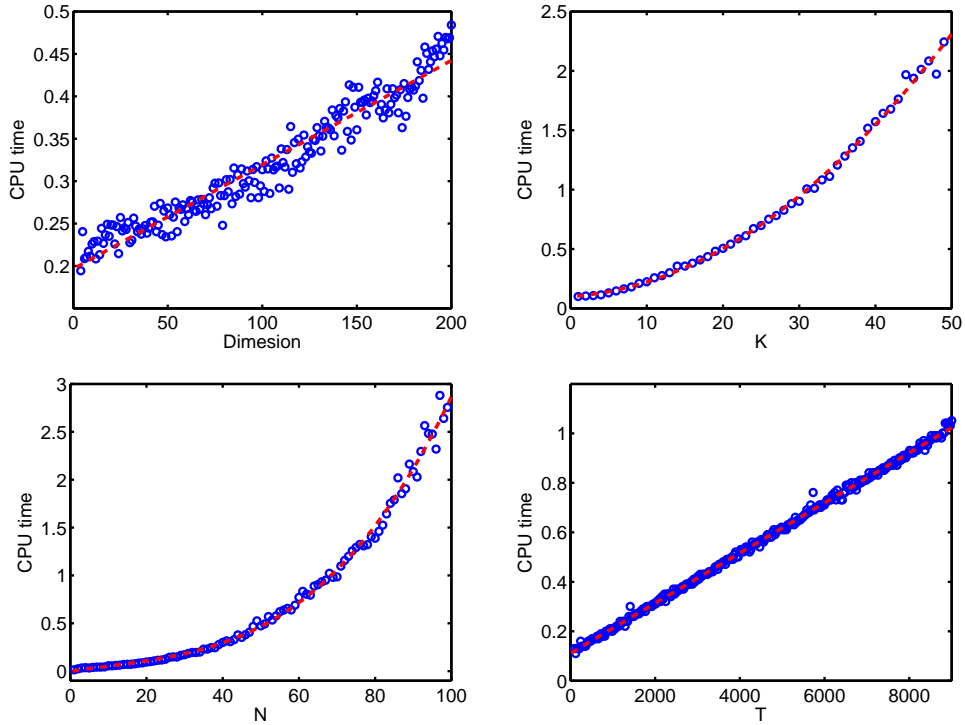
**Fig. 8** Influence of the *regularization parameter*  $\epsilon^2$  on the identified *cluster affiliation function*  $\gamma^\epsilon$ . Grayscale represents the values of the function  $\gamma^\epsilon(t)$  calculated for different values of  $\epsilon^2$  at different times calculated for the time series generated with  $\mathbf{K} = 2, T = 1250, N = 50, m = 1, \alpha = 25$  and cluster switching from Fig. 6. Dashed lines denote the moments when the original process from Fig. 6 was switching between the clusters.

This will obviously prohibit the application of the approaches like GMMs or HMMs (implying a priori probabilistic assumptions on the data) to time series analysis of  $x_t$  in considered case.

Fig. 3 demonstrates the cluster affiliations for the time series  $x_t$  (lower panel of Fig. 2) calculated with two standard clustering methods: K-means [10, 11] and fuzzy C-Means algorithms [12, 11]. As can be seen from the comparison of the upper panel of Fig. 2 and Fig. 3, both of the standard methods



**Fig. 9**  $l_2$  norm of the difference between the original switching process from Fig. 6 and its estimate based on data analysis for different angles  $\alpha$  and different ratios of dominant covariance matrix eigenvalues  $\lambda_1$  and  $\lambda_2$ .



**Fig. 10** Statistics of the computational performance obtained from 1000 realizations of the model process with  $\alpha = 25$ ,  $\epsilon^2 = 0.1$  generated for a switching process according to Fig. 6. Dashed lines represent the polynomials of respective order fitted to the simulation data (circles). Four panels demonstrate the performance of the method: (i) wrt. observation dimension  $n$  for  $\mathbf{K} = 2$ ,  $m = 1$ ,  $T = 1250$ ,  $N = 50$  ( $\mathcal{O}(n)$ , upper left panel), (ii) wrt. number of clusters  $\mathbf{K}$  for  $m = 1$ ,  $T = 1250$ ,  $N = 50$ ,  $n = 100$  ( $\mathcal{O}(\mathbf{K}^2)$ , upper right panel), (iii) wrt. number of finite elements  $N$  for  $\mathbf{K} = 2$ ,  $m = 1$ ,  $T = 1250$ ,  $n = 100$  ( $\mathcal{O}(N^2 \log(N))$ , lower left panel), and (iv) wrt. length  $T$  of the time series for  $\mathbf{K} = 2$ ,  $m = 1$ ,  $N = 50$ ,  $n = 100$  ( $\mathcal{O}(T)$ , lower right panel).

fail to recover the original hidden process used in the generation of the time series  $x_t$ . This is due to the strong overlap between the observed processes in both states. Fig. 4 demonstrates the results of the FEM-based minimization of the optimization problem (19,20, 21) for  $N = T = 500$ ,  $\mathbf{K} = 2$ , with model distance functional  $g(x_t, \theta_i)$  of the form (6) ( $\theta_i \in \mathbf{R}^1$ ,  $i = 1, 2$  are one-dimensional Euclidean centers of the respective cluster states) and different values of the regularization parameter  $\epsilon^2$ . In contrast to the

standard clustering techniques (imposing no requirements on the persistence of the process  $\Gamma(t)$ ), FEM-clustering procedure described above allows to recover the correct form of the hidden transition process. Fig. 4 also demonstrates the effect of regularization on the identified transition process: increasing smoothness and persistence of  $\Gamma(t)$  with growing  $\epsilon^2$ .

Next, we consider a problem of optimal  $\mathbf{K}$  identification. As was mentioned in Sec. 4, in order to determine the upper bound  $\mathbf{K}_{max}$  of the *number of statistically distinguishable* cluster states for any fixed  $\epsilon^2$ , five following steps are needed: (i) some (high) number of cluster states  $\mathbf{K}$  should be chosen as a start value a priori, (ii) FEM-clustering procedure should be performed, (iii) a posteriori probabilistic assumptions about the observation data  $x_t$  in each of the  $\mathbf{K}$  identified cluster states should be stated, (iv) confidence intervals for the identified parameters  $\theta_i$  in each of the  $\mathbf{K}$  identified cluster states should be calculated and (v) in the case of their overlap in all components for a certain pair of the identified states, the number of cluster states should be set to  $\mathbf{K} = \mathbf{K} - 1$  and the whole procedure should be repeated from the step (ii). Since we choose the *model distance functional*  $g(x_t, \theta_i)$  of the form (6), parameters  $\theta_i$  are calculated in the standard way of K-Means procedure, as *conditional expectation values* of the observed process in each of the states). Standard tools of statistical hypothesis testing under i.i.d.-assumption of the observed data in each of the clusters can be applied to estimate the confidence intervals for each  $\theta_i, i = 1, \dots, \mathbf{K}$  and to test the statistical distinguishability of the cluster sets resulting from the FEM-clustering procedure a posteriori [39]. Left panel of Fig. 5 demonstrates the outcome of the described procedure: it shows that for a wide range of the regularization parameters  $\epsilon^2$ , there are at most two statistically distinguishable cluster states. It coincides with the number of the local models used in the generation of the original data  $x_t$ .

Right panel of Fig. 5 demonstrates the dependence between the regularization parameter  $\epsilon^2$ , respective value of the averaged clustering functional (10) at the optimum of the regularized problem (19,20, 21) and the mean exit time of identified states. As can be seen from Fig. 5, shape of the resulting dependence between the quantities is described by the L-curve in 3 dimensions. Analogously with the case of Tikhonov-regularized linear problems, optimal regularization parameter ( $\epsilon^2 = 0.2$  for the considered example) can be determined from the edge-point of the resulting L-curve.

Finally, we compare the numerical cost of the presented FEM-clustering method with standard clustering techniques (like K-Means and fuzzy C-Means). As is known from the literature [10, 12, 11], both K-means and fuzzy C-Means algorithms scale as  $\mathbf{O}(T)$  (where  $T$  is the length of the time series). Careful inspection of the regularized problem (19,20, 21) demonstrates that the limiting step of the described FEM-based algorithmic procedure is solution of the sparse quadratic minimization problem with constraints, scaling as  $\mathbf{O}(N \log(N))$  (where  $N$  is the number of finite element functions,  $N \leq T$ ). Comparison of scalings reveals that the FEM-clustering procedure will be computationally more efficient than standard clustering approaches in two cases: (i) if the approximation threshold  $\delta_N$  is set to be large (dependent on the problem resulting in a low number of finite element functions  $N$ ) or (ii) if the underlying process is *persistent* (implying that only few finite element functions are needed to resolve the hidden process  $\Gamma(t)$ ). In all other cases presented FEM-clustering method is numerically more expensive. In the considered test example (with  $N$  being explicitly set to be equal to the number of time steps in the analyzed time series,  $N = T = 500$ ), FEM-clustering procedure is twice as slow as compared to the K-Means method. However, as was already mentioned above, FEM-clustering procedure revealed the correct persistent dynamics of the hidden process  $Y_t$ , whereas both of the standard procedures fail (see Fig. 3 and Fig. 4).

## 5.2 Multi-dimensional test data: computational performance

Now we will apply the method to a test system consisting of a given discrete process switching between two  $n$ -dimensional data clusters (see Fig. 6). Data in clusters is distributed according to a tensor product of a two-dimensional Gaussian distribution (7) with a  $(n - 2)$ -dimensional uniform distribution on the interval  $[0, 10]$ . This results in a multivariate *non-Gaussian distribution* in  $n$  dimensions for each of the

cluster states. Both distributions have the same expectation value

$$\mu_{1,2} = \begin{pmatrix} 10 \\ 3 \\ 5 \\ 5 \\ \dots \\ 5 \end{pmatrix}, \quad (34)$$

and one of the covariance matrices (in both Gaussian degrees of freedom) is a rotation of the other covariance matrix through some predefined angle  $\alpha$ :

$$\Sigma_1 = \begin{pmatrix} 2 & 0.5 \\ 0.5 & 1 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} \cos(\alpha) & \sin(\alpha) \\ -\sin(\alpha) & \cos(\alpha) \end{pmatrix} \Sigma_1, \quad (35)$$

Covariances in  $(n - 2)$  non-Gaussian degrees of freedom are identical for both distributions.

The resulting  $n$ -dimensional non-Gaussian distributions are chosen in such a way that neither the *geometric clustering* approaches (like K-Means or fuzzy C-Means) nor the *GMMs/HMMs* will be able to cluster the data properly. Such kinds of distributions are very common in computational finance, e.g., in the analysis of stock returns where the change of the market phase is connected to the change of volatility of the underlying stochastic process [6]. Therefore in the following we will use the *EOF model distance functional* (9) (with  $\theta_i$  being the projectors on dominant one dimensional linear manifolds  $\mathcal{T}_i \in \mathbf{R}^{n \times m}, m = 1$ ) with *linear finite elements* (27) for the clustering of the resulting time series. As an example, Fig. 7 demonstrates two time series in two Gaussian degrees of freedom generated with the help of the switching process from Fig. 6 for two different values of the angle  $\alpha$  (90 and 25). All other degrees of freedom are generated according to the same uniform distribution and look non-distinguishable.

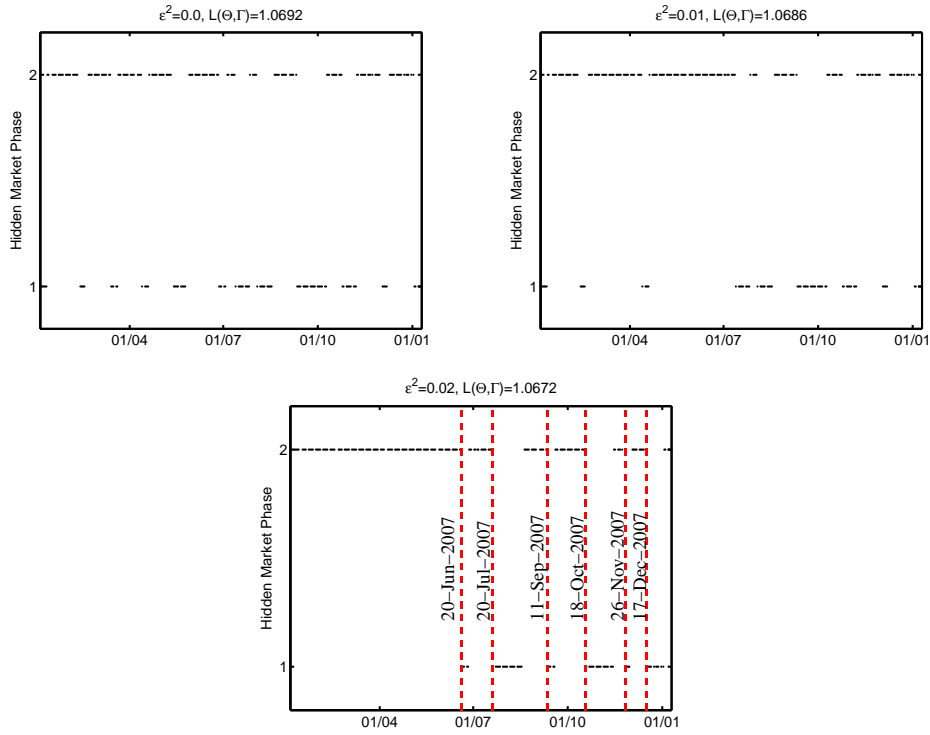
Fig. 8 demonstrates the effect of regularization on the optimization process: as was already discussed above, an increase of  $\epsilon^2$  leads to a growing metastability of the resulting *cluster affiliation function*. As we see from Fig. 8 this results in a *coarse graining of the identified affiliation functions*, i. e., only "long living" structures in  $\gamma$  "survive" with increasing  $\epsilon^2$ . Fig. 9 demonstrates the sensitivity of the optimization procedure to the input time series, i. e., it shows the clustering error as a function of two variables: (i) angle  $\alpha$  between the clusters and (ii) ratio of two dominant eigenvalues. As expected, the quality of clustering is increasing with increasing  $\alpha$  and decreasing  $\lambda_2/\lambda_1$ . This is explained by the growing numerical separability of the dominant subspaces in the context of the *EOF model distance functional*.

Fig. 10 demonstrates the computational performance of the resulting clustering method measured experimentally from 1.000 different realizations of the analyzed model trajectories (all generated with the same cluster affiliation function, see Fig. 6): it is linear in the observation dimension and time series length, quadratic in the number of clusters and polynomial in the number  $N$  of finite elements (scales approx.  $\mathbf{O}(N^2 \log(N))$ ). It should be mentioned that the standard QP-solver from MATLAB was applied in the current realizations of the code. The sparsity structure of the QP subproblem allows to use *sparse QP (SQP)* solvers available on the market. This will reduce the numerical cost of the method to  $\mathbf{O}(N \log(N))$ .

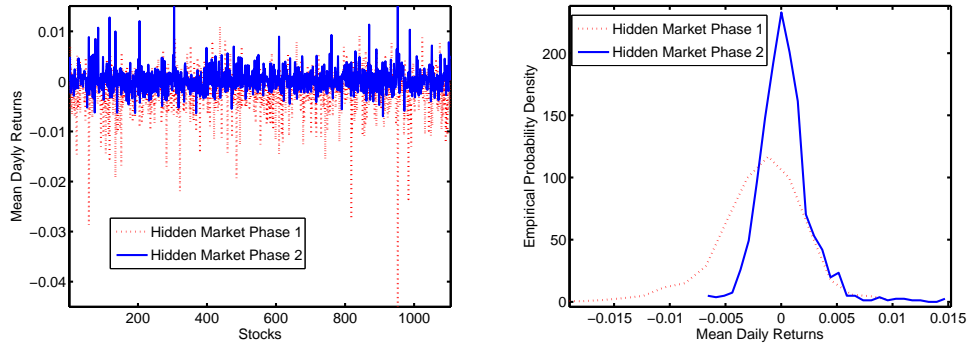
### 5.3 Analysis of stock daily returns: NASDAQ

Finally we will apply the numerical method to identify the hidden market phases based on daily returns of the 1106 stocks from the *NASDAQ* stock exchange between Jan.03 2007 and Jan.10 2008 (data is acquired from <http://finance.yahoo.com>). Our aim is to identify the *hidden market phases* and to interpret them in the context of global market dynamics. Due to the fact that the time series has only 257 elements (because there are 257 trading days in the analyzed time interval), we do not expect to find more than few statistically distinguishable clusters. We apply the *EOF model distance functional* (9) (with  $\theta_i$  being the projectors on dominant one dimensional linear manifolds  $\mathcal{T}_i \in \mathbf{R}^{1106 \times 1}$ ) and *linear finite elements* (27) with  $\mathbf{K} = 2, m = 1, T = 257, N = 50, n = 1106$ . Because of the fact that the proposed numerical framework approaches only towards a *local minimum* of the *regularized averaged*





**Fig. 11** Hidden market phases for different values of the *regularization parameter*  $\epsilon^2$  as calculated from historical time series of daily NASDAQ stock returns using the *EOF model distance functional* (9) ( $\mathbf{K} = 2, n = 1106, N = 50, m = 1, T = 257$ , in each case the optimization was repeated 100 times and the solution with the lowest value of *averaged clustering functional*  $\mathbf{L}(\Theta, \Gamma)$  is taken in each case).



**Fig. 12** Left panel: mean daily returns for 1106 NASDAQ titles calculated as *conditional averages* over the corresponding hidden market phases from the lower panel of Fig. 11. Right panel: empirical probability distributions of mean daily returns in both identified market phases.

*clustering functional* (dependent on the initial parameter values chosen for optimization), we repeat the optimization 100 times with different randomly generated initial guesses for the parameters and keep the result with the lowest value of the *averaged clustering functional* (10).

Results of this procedure for different values of  $\epsilon^2$  are demonstrated in Fig. 11. As can be seen, the lowest value of the *averaged clustering functional* is achieved for the optimization with highest *regularity*  $\epsilon^2 = 0.02$ , this also means the highest *metastability* of the process switching between the market phases. Further increase of  $\epsilon^2$  leads to a complete domination of the regularization part in optimization and

suppression of the part corresponding to the *averaged clustering functional*. This results in identification of the hidden path with no transitions at all and higher values of (10). Therefore we can assume that the identified hidden path for  $\epsilon^2 = 0.02$  is optimal wrt. the *averaged clustering functional* (10), i. e., it has a lowest value of  $\mathbf{L}$  among all other identified pathways.

To interpret the resulting hidden path in terms of the global market dynamics we will first have a look at the *mean daily stock returns* [6] for each of the phases. The right panel of Fig. 12 demonstrates that the empirical probability distribution in the second market phase is narrow with a "heavy tail" in the positive direction, whereas its counterpart for the first hidden state is much wider with a heavy tail in the negative direction. This means that the second market phase corresponds to a more stable, positive global dynamics and the first phase is characterized by much more unstable, volatile and negative dynamics. Inspection of the switches between the market phases reveals that the first transition to a market phase 1 happens on 20-Jun-2007 which approximately corresponds to the beginning of the US subprime mortgage financial crisis among the US hedge funds.

## 6 Conclusion

We have presented a numerical framework for the clustering of multidimensional time series based on minimization of a *regularized averaged clustering functional*. Finite element discretization of the problem allowed us to suggest a numerical algorithm based on the splitting procedure applicable for a wide class of clustering problems. We have investigated the conditions under which the proposed numerical method is monotone and analyzed the connection between the *regularization factor* and *metastability* in context of homogeneous Markov-jump processes.

One of the open problems is a rigorous mathematical investigation of the discretization error. It is appealing to apply the asymptotical theory developed for partial differential equations. This will allow to construct much more efficient *adaptive numerical methods* of data-clustering. Another problem is the *locality* of the proposed numerical scheme, i. e., the obtained result is dependent on the initial value. In the current implementation we solve the problem by "brute force", just repeating the optimization many times with different randomly initialized parameter values. Finally, the sparsity of the matrices involved in QP-subproblem allows to use much more efficient *sparse quadratic programming (SQP)* tools which are very suitable for parallel processing, this is also a matter of future research.

Working with multidimensional data, it is very important to be able to extract some reduced description out of it (e.g., in form of *essential degrees of freedom* or *hidden pathes*). In order to control the reliability of obtained results, one has to analyze the sensitivity of results wrt. the length of the time series and the number  $K$  of the hidden states. We have given some hints for the selection of an optimal  $K$  and explained how the quality of the resulting reduced representation can be acquired.

## Acknowledgement

The author thanks O. Sander (FU Berlin) and two unknown referees for careful reading of the manuscript and for helpful discussion. The work was supported by the DFG research center MATHEON "Mathematics for key technologies" in Berlin and SPP "MetStroem".

## References

- [1] Sutera A. Benzi R., Parisi G. and Vulpiani A. Stochastic resonance in climatic change. *Tellus*, 3:10–16, 1982.
- [2] C. Nicolis. Stochastic aspects of climatic transitions-response to a periodic forcing. *Tellus*, 34:1–+, 1982.
- [3] A.A. Tsonis and J.B. Elsner. Multiple attractors, fractal basins and longterm climate dynamics. *Beit. Phys. Atmos.*, 63:171–176, 1990.
- [4] T.N. Palmer. A Nonlinear Dynamical Perspective on Climate Prediction. *Journal of Climate*, 12:575–591, February 1999.
- [5] J.D. Hamilton. A new approach to the economic analysis of nonstationary time series and the business cycle. *Econometrica*, 57:357–384, 1989.
- [6] R.S. Tsay. *Analysis of financial time series*. Wiley-Interscience, 2005.

- [7] R. Elber and M. Karplus. Multiple conformational states of proteins: A molecular dynamics analysis of Myoglobin. *Science*, 235:318–321, 1987.
- [8] H. Frauenfelder, P. J. Steinbach, and R. D. Young. Conformational relaxation in proteins. *Chem. Soc.*, 29A:145–150, 1989.
- [9] Ch. Schütte, A. Fischer, W. Huisinga, and P. Deuffhard. A direct approach to conformational dynamics based on hybrid Monte Carlo. *J. Comput. Phys.*, 151:146–168, 1999.
- [10] J.C. Bezdek. *Pattern recognition with fuzzy objective function algorithms*. Plenum Press, New York, 1981.
- [11] F. Höppner, F. Klawonn, R. Kruse, and T. Runkler. *Fuzzy cluster analysis*. John Wiley and Sons, New York, 1999.
- [12] J.C. Bezdek, R.H. Hathaway, M.J. Sabin, and W.T. Tucker. Convergence theory for fuzzy c-mens: counterexamples and repairs. *IEEE Trans. Systems.*, 17:873–877, 1987.
- [13] R.H. Hathaway and J.C. Bezdek. Switching regression models and fuzzy clustering. *IEEE Transactions on Fuzzy Systems.*, 1:195–204, 1993.
- [14] G.J. McLachlan and D. Peel. *Finite Mixture Models*. Wiley, 2000.
- [15] S. Frühwirth-Schnatter. *Finite Mixture and Markov Switching Models*. Springer, 2006.
- [16] J.A. Bilmes. *A Gentle Tutorial of the EM Algorithm and its Applications to Parameter Estimation for Gaussian Mixture and Hidden Markov Models. Thechnical Report*. International Computer Science Institute, Berkeley, 1998.
- [17] I. Horenko, J. Schmidt-Ehrenberg, and Ch. Schütte. Set-oriented dimension reduction: Localizing Principal Component Analysis via Hidden Markov Models. In R. Glen M.R. Berthold and I. Fischer, editors, *CompLife 2006*, volume 4216 of *Lecture Notes in Bioinformatics*, pages 98–115. Springer, Berlin Heidelberg, 2006.
- [18] I. Horenko, R. Klein, S. Dolaphtchiev, and Ch. Schuette. Automated generation of reduced stochastic weather models i: simultaneous dimension and model reduction for time series analysis. *SIAM Mult. Mod. Sim.*, 6(4):1125–1145, 2008.
- [19] I. Horenko. On simultaneous data-based dimension reduction and hidden phase identification. *To appear in J. Atmos. Sci.*, 2008. (available via biocomputing.mi.fu-berlin.de).
- [20] A.H. Monahan. Nonlinear principal component analysis by neural networks: Theory and application to the lorenz system. *J. Climate*, 13:821–835, 2000.
- [21] W. Härdle and L. Simar. *Applied Multivariate Statistical Analysis*. Springer, New York, 2003.
- [22] I.T. Jolliffe. *Principal Component Analysis*. Springer, 2002.
- [23] P. Deuffhard. *Newton Methods for Nonlinear Problems. Affine Invariance and Adaptive Algorithms*, volume 35 of *Computational Mathematics*. Springer, Heidelberg, 2004.
- [24] A. Tikhonov. On the stability of inverse problems. *Dokl. Akad. Nauk SSSR*, 39(5):195–198, 1943.
- [25] A. Hoerl. Application of ridge analysis to regression problems. *Chemical Engineering Progress*, 58:54–59, 1962.
- [26] G. Wahba. *Spline Models for Observational Data*. SIAM, 1990.
- [27] D. Braess. *Finite Elements: Theory, Fast Solvers and Applications to Solid Mechanics*. 3rd. edition.
- [28] M.K. Kozlov and S.P. Tarasov L.G. Kachiyan. Polynomial solvability of convex quadratic programming. *Sov. Math. Dokl.*, 20:1108–1111, 1979.
- [29] J. Nocedal and S.J. Wright. *Numerical Optimization*. Springer, 1999.
- [30] B. de Finetti. *Theory of Probability. Vols. I and II*. Wiley, New York, 1974.
- [31] H. Gardiner. *Handbook of stochastic methods*. Springer, Berlin, 2000.
- [32] D. Calvetti, S. Morigi, L. Reichel, and F. Sgallari. Tikhonov regularization and the l-curve for large discrete ill-posed problems. *J. Comput. Appl. Math.*, 123(1-2):423–446, 2000.
- [33] D. Calvetti, L. Reichel, and A. Shuibi. L-Curve and Curvature Bounds for Tikhonov Regularization. *Numerical Algorithms*, 35:301–314, 2004.
- [34] P. Metzner, I. Horenko, and Ch. Schuette. Generator estimation of Markov Jump processes based on incomplete observations nonequidistant in time. *Phys. Rev. E*, 76:0667021, 2007.
- [35] Christof Schütte and Wilhelm Huisinga. Biomolecular conformations can be identified as metastable sets of molecular dynamics. In P. G. Ciarlet and J.-L. Lions, editors, *Handbook of Numerical Analysis*, volume X, pages 699–744. Elsevier, 2003.
- [36] M. Weber and P. Deuffhard. Perron cluster cluster analysis. *J. Chem. Phys.*, 5:802–827, 2003.
- [37] H. Engl and A. Neubauer. An improved version of Marti’s method for solving ill-posed linear integral equations. *Math. Comp.*, 45:405–416, 1985.
- [38] Eberhard Schock. Morozovs discrepancy principle for tikhonov regularization of severely ill-posed problems in finite-dimensional subspaces. *Numer. Funct. Anal. Optim.*, 21:901–916, 2000.
- [39] B. Kedem and K. Fokianos. *Regression models for time series analysis*. Wiley Series in Probability and Statistics, 2002.