# State-space dynamics distance for clustering sequential data

Darío García-García *, Emilio Parrado-Hernández, Fernando Diaz-de-Maria

Signal Theory and Communications Deparment, Escuela Politécnica Superior, Universidad Carlos III de Madrid, Avda. de la Universidad, 30, 28911 Leganés, Spain

## A R T I C L E   I N F O

## A B S T R A C T

This paper proposes a novel similarity measure for clustering sequential data. We first construct a common state space by training a single probabilistic model with all the sequences in order to get a unified representation for the dataset. Then, distances are obtained attending to the transition matrices induced by each sequence in that state space. This approach solves some of the usual overfitting and scalability issues of the existing semi-parametric techniques that rely on training a model for each sequence. Empirical studies on both synthetic and real-world datasets illustrate the advantages of the proposed similarity measure for clustering sequences.

© 2010 Elsevier Ltd. All rights reserved.

## 1. Introduction

Clustering is a core task in machine learning and data processing. It is an unsupervised technique whose goal is to unveil a natural partition of data into a number of groups or *clusters*. Clustering techniques have been widely studied and applied for a long time, yielding a large number of well-known and efficient algorithms [1]. Recently, the family of algorithms collectively known as spectral clustering [2,3] has received a lot of attention due to its good performance and solid theoretical foundation. These methods share a graph-theoretic approach that results in non-parametric partitions of a dataset, in the sense that they do not require any parametric model of the input data.

The application of clustering techniques to sequential data is lately receiving growing attention [4]. In these scenarios, the useful information is not encoded only in the data vectors themselves, but also in the way they evolve along a certain dimension (usually time). Examples of sequential data range from stock market analysis to audio signals, video sequences, etc. Developing machine learning techniques for these scenarios poses additional difficulties compared to the classical setting where data vectors are assumed to be independent and identically distributed (i.i.d.). The common framework consists in combining a distance or similarity measure between sequences of variable length, which captures their dynamical behaviour, with a clustering algorithm developed for i.i.d. data.

The design of similarity measures for sequences is generally addressed from a model-based perspective. Specifically, hidden Markov models (HMMs) [5] are usually employed as models for the sequences in the dataset. HMMs have been widely used in signal processing and pattern recognition because they offer a good trade-off between complexity and expressive power. Following the work of Smyth [6], many researchers [7–10] have proposed different distance measures based on a likelihood matrix obtained using each single sequence to train an HMM. Besides, [11] defines the similarity between two sequences as the probability product kernel (PPK) between the HMMs trained on each sequence. The application of a non-parametric clustering to the distance matrix obtained in one of the aforementioned ways yields a semi-parametric model: each individual sequence follows a parametric model, but no assumption is made about the global cluster structure. These semi-parametric methods have been shown [10,11] to outperform both fully parametric methods such as mixture of HMMs [12] or combinations of HMMs and dynamic time warping [13].

This paper aims at solving a main weakness of the aforementioned semi-parametric models. The fact that each model is trained using just one sequence can lead to severe overfitting or non-representative models for short or noisy sequences. In addition, the learning of a semi-parametric model involves the calculation of a number of likelihoods or probability product kernels that is quadratic in the number of sequences, which hinders the scalability of the method. To overcome these disadvantages, we propose to train a single HMM using all the sequences in the dataset, and then cluster the sequences attending to the transition matrices that they induce in the state-space of the common HMM.

This approach is radically different from the aforementioned methods (including the authors' previous work [7]) in the sense that it is not based on likelihoods, but on divergences between the transition probabilities that each sequence induces under the common model. In other words, we no longer evaluate the likelihoods of the sequences on some models and then define the distance accordingly. Instead, the focus is now shifted towards

* Corresponding author. Tel.: +34 91 6248805; fax: +34 91 6248749.
E-mail addresses: dggarcia@tsc.uc3m.es (D. García-García), emipar@tsc.uc3m.es (E. Parrado-Hernández), fdiaz@tsc.uc3m.es (F. Diaz-de-Maria).
URL: http://www.tsc.uc3m.es/~dggarcia (D. García-García).

parameter space. Moreover, the identification of each sequence with a transition matrix opens up new possibilities since the metric can be based on the short term transitions, the long term stationary state distribution or on some middle ground.

The rest of the paper is organized as follows: Section 2 starts with a brief review of hidden Markov models and spectral clustering followed by a presentation of the state of the art in model-based clustering of sequences. Section 3 describes how to cluster sequences using information about their dynamics in a common state space. This new proposal is empirically evaluated in comparison with other methods in Section 4. Finally, Section 5 draws the main conclusions of this work and sketches some promising lines for future research.

## 2. Clustering sequences with hidden Markov models

This section reviews the state-of-the-art framework employed to cluster sequential data, which consists of two phases: (1) the design of a similarity or distance measure between sequences based on hidden Markov models; and (2) the use of that measure in a clustering algorithm. We opt for spectral clustering due to its good results, as reported in the literature. Nonetheless, the distances presented in this work can be used in combination with any clustering algorithm.

### 2.1. Hidden Markov models

Hidden Markov models (HMMs) [5] are a type of parametric, discrete state-space model. They provide a convenient model for many real-life phenomena, while allowing for low-complexity algorithms for inference and learning. Their main assumptions are the independence of the observations given the hidden states and that these states follow a Markov chain.

Assume a sequence $\mathbf{S}$ of $T$ observation vectors $\mathbf{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$. The HMM assumes that $\mathbf{x}_t$, the $t$th observation of the sequence, is generated according to the conditional emission density $p(\mathbf{x}_t|q_t)$, with $q_t$ being the hidden state at time $t$. The state $q_t$ can take values from a discrete set $\{s_1, \dots, s_K\}$ of size $K$. The hidden states evolve following a time-homogeneous first-order Markov chain, so that $p(q_t|q_{t-1}, q_{t-2}, \dots, q_0) = p(q_t|q_{t-1})$.

In this manner, the parameter set $\theta$ that defines an HMM consists of the following distributions:

- The initial probabilities vector $\boldsymbol{\pi} = \{\pi_i\}_{i=1}^K$, where $\pi_i = p(q_0 = s_i)$.
- The state transition probability, encoded in a matrix $\mathbf{A} = \{a_{ij}\}_{i,j=1}^K$ with $a_{ij} = p(q_{t+1} = s_j|q_t = s_i)$, $1 \le i, j \le K$.
- The emission pdf for each hidden state $p(\mathbf{x}_t|q_t = s_i)$, $1 \le i \le K$.

From these definitions, the likelihood of a sequence $\mathbf{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ can be written in the following factorized way:

$$p(\mathbf{S}|\theta) = \sum_{q_0, \dots, q_T} \pi_{q_0} p(\mathbf{x}_0|q_0) \prod_{t=1}^{T} p(\mathbf{x}_t|q_t) a_{q_{t-1}, q_t}. \tag{1}$$

The training of this kind of models in a maximum likelihood setting is usually accomplished using the Baum–Welch method [5], which is a particularization of the well-known EM algorithm. The E-step finds the expected state occupancy and transition probabilities, which can be done efficiently using the forward–backward algorithm [5]. This algorithm implies the calculation of both the forward $\alpha$ and backward $\beta$ variables that are defined as follows:

$$\alpha_k(t) = p(\mathbf{x}_1, \dots, \mathbf{x}_t, q_t = s_k), \tag{2}$$

$$\beta_k(t) = p(\mathbf{x}_{t+1}, \dots, \mathbf{x}_T|q_t = s_k). \tag{3}$$

These variables can be obtained in $O(K^2 T)$ time through a recursive procedure and can be used to rewrite the likelihood from Eq. (1) in the

following manner:

$$p(\mathbf{S}|\theta) = \sum_{k=1}^{K} \alpha_k(t)\beta_k(t), \tag{4}$$

which holds for all values of $t \in \{1, \dots, T\}$.

Given a previously estimated $\mathbf{A}$, the state transition probabilities can be updated using the forward/backward variables and that previous estimation, yielding

$$\tilde{a}_{ij} \propto \sum_{t'=1}^{T} \alpha_i(t') a_{ij} p(\mathbf{x}_{t'+1}|q_{t'+1} = s_j)\beta_j(t'+1). \tag{5}$$

Then, the M-step updates the parameters in order to maximize the likelihood given the expected hidden states sequence. These two steps are then iterated until convergence. It is worth noting that the likelihood function can have many local maxima, and this algorithm does not guarantee convergence to the global optimum. Due to this, it is common practice to repeat the training several times using different initializations and then select as the correct run the one providing a larger likelihood.

The extension of this training procedure to multiple input sequences is straightforward. The interested reader is referred to [5] for a complete description.

### 2.2. Spectral clustering

Clustering [1] consists in partitioning a dataset $\mathcal{S}$ comprised $N$ elements into $C$ disjoint groups called clusters. Data assigned to the same cluster must be similar and, at the same time, distinct from data assigned to the rest of clusters. It is an unsupervised learning problem, meaning that it does not require any prior labelling of the data, and thus it is very appropriate for exploratory data analysis or scenarios where obtaining such a labelling is costly.

Algebraically, a clustering problem can be formulated in the following way. Given a dataset $\mathcal{S}$, one forms a $N \times N$ similarity matrix $\mathbf{W}$, whose $ij$th element $w_{ij}$ represents the similarity between the $i$th and $j$th instances. The clustering problem then consists in obtaining a $N \times C$ clustering matrix $\mathbf{Z}$, where $z_{ic} = 1$ if instance $i$ belongs to cluster $c$ and $z_{ic} = 0$ otherwise, which is optimal under some criteria.

Spectral clustering (SC) algorithms [3] approach the clustering task from a graph-theoretic perspective. Data instances form the nodes $V$ of a weighted graph $G = (V,E)$ whose edges $E$ represent the similarity or adjacency between data, defined by the matrix $\mathbf{W}$. This way, the clustering problem is cast into a graph partitioning one. The clusters are given by the partition of $G$ into $C$ groups that optimize certain criteria such as the normalized cut [14]. Finding such an optimal partition is an NP-hard problem, but it can be relaxed into an eigenvalue problem on the Laplacian matrix $\mathbf{L} = \mathbf{D} - \mathbf{W}$, where $\mathbf{D}$ is the diagonal matrix with elements $d_{ii} = \sum_{j=1}^{N} w_{ij}$, or one of its normalized versions, followed by $k$-means [20] or any other clustering algorithm on the rows of the matrix of selected eigenvectors. Specifically, in the experiments in this paper, we employ the spectral clustering algorithm defined in [2], which uses the symmetric version of the normalized Laplacian $\mathbf{L}_{sym} = \mathbf{D}^{-1/2}\mathbf{L}\mathbf{D}^{-1/2}$. The actual clustering is applied to the normalized rows of the $C$ eigenvectors associated with the lowest eigenvalues. It is worth noting that the embedding of the original data into the rows of the normalized Laplacian eigenvectors also have some appealing properties from the point of view of dimensionality reduction [15].

The time complexity for the spectral clustering is dominated by the eigendecomposition of the normalized Laplacian, which in general is $O(N^3)$. However, if the affinity matrix is sparse (e.g. if only the affinities between the nearest neighbours of a given node are

considered), there exist efficient iterative methods that notably reduce this complexity, such as the Lanczos method [16], which makes it feasible even for large datasets [3].

### 2.3. Semi-parametric model-based sequence clustering

Semi-parametric sequence clustering methods make no assumptions about the cluster structure. These methods typically use generative models such as HMMs in order to represent the sequences in the dataset in a common space, and then apply a non-parametric clustering algorithm. Following [6], many methods use a log-likelihood matrix $\mathbf{L}$ whose $ij$th element $l_{ij}$ is defined as

$$l_{ij} = \log p_{ij} = \frac{1}{\text{length}(\mathbf{S}_j)} \log p(\mathbf{S}_j|\theta_i), \quad 1 \leq i, \ j \leq N, \tag{6}$$

where $\theta_i$ is the model trained for the $i$th sequence $\mathbf{S}_i$ and $N$ is the number of sequences in the dataset. The log-likelihoods are normalized with respect to the length of the sequences to compensate the exponentially increasing size of the probability space. Recall that the cost of such an individual likelihood calculation for a sequence $\mathbf{S}_i$ is $O(K^2 T_i)$, where $T_i$ is the length of $\mathbf{S}_i$.

The literature reports several methods to construct a similarity (distance) matrix $\mathbf{D}$ from $\mathbf{L}$, such as:

- The "symmetrized distance" (SYM) [6]: $d_{ij}^{\text{SYM}} = \frac{1}{2}(l_{ij} + l_{ji})$.
- The BP distance in [8]: $d_{ij}^{\text{BP}} = \frac{1}{2}\{(l_{ij} - l_{ii})/l_{ii} + (l_{ji} - l_{jj})/l_{jj}\}$, which takes into account how well a model represents the sequence it has been trained on.
- The Yin–Yang distance of [10] $d_{ij}^{\text{YY}} = |l_{ii} + l_{jj} - l_{ij} - l_{ji}|$.

The distance matrix $\mathbf{D}$ is then fed into a clustering algorithm that partitions the set of sequences. It is worth noting that all these methods define the distance between two sequences $\mathbf{S}_i$ and $\mathbf{S}_j$ using solely the models trained on these particular sequences ($\theta_i$ and $\theta_j$).

In [11] another method for constructing the similarity matrix in a model-based approach is proposed which avoids the calculation of the likelihood matrix $\mathbf{L}$. Again an HMM is trained on each individual sequence, but then similarities between sequences are computed directly through a probability product kernel (PPK). Specifically, this kernel is obtained by integrating a product of the probability distributions spanned by two HMMs in the space of the sequences of a fixed length $T_{\text{PPK}}$, which is a free parameter. This way, $\mathbf{L}$ is no longer necessary because the similarities are obtained directly from parameters of the models. The calculation of the PPK between two HMMs can be carried out in $O(K^2 T_{\text{PPK}})$ time.

Furthermore, the method proposed in [7] assumes the existence of a latent model space $\theta$ formed by some HMMs that actually span the data space. Then, the models $\theta_1, \ldots, \theta_N$ trained on each sequence are regarded as a set $\tilde{\theta}$ of intelligently sampled points of $\theta$. Let $\mathbf{L}_N$ be a column-wise normalization of $\mathbf{L}$. The $i$th column of $\mathbf{L}_N$ can be interpreted as a probability density function $f_{\tilde{\theta}}^{\mathbf{S}_i}(\theta)$ over the approximated model space $\tilde{\theta}$ for $\mathbf{S}_i$. This way it is possible to re-express $\mathbf{L}_N$ as

$$\mathbf{L}_N = [f_{\tilde{\theta}}^{\mathbf{S}_1}(\theta), \ldots f_{\tilde{\theta}}^{\mathbf{S}_N}(\theta)].$$

This interpretation leads to a distance measure consisting in Kullback–Leibler (KL) divergences between the columns of $\mathbf{L}_N$, so

$$d_{ij}^{\text{KL}} = D_{\text{KL}_{\text{SYM}}}(f_{\tilde{\theta}}^{\mathbf{S}_i}||f_{\tilde{\theta}}^{\mathbf{S}_j}), \tag{7}$$

where $D_{\text{KL}_{\text{SYM}}}$ stands for the symmetrized version of the KL divergence: $D_{\text{KL}_{\text{SYM}}}(p||q) = \frac{1}{2}(D_{\text{KL}}(p||q) + D_{\text{KL}}(q||p))$. This way, the distance between two sequences is obtained in a global way, using a data-dependent representation. Under this paradigm, it is possible to select a subset of $P \leq N$ sequences to train individual models on, instead of fitting an HMM to every sequence in the dataset. This can reduce the

computational load and improve the performance on real-world data. However, the a priori selection of such a subset remains an open problem.

The aforementioned semi-parametric methods share the need to train an individual HMM on each sequence in the dataset (or in a subset of it, as in [7]). We see this as a disadvantage for several reasons. Short sequences are likely to lead to overfitted models, providing unrealistic results when used to evaluate likelihoods or PPKs. Moreover, training individual models in an isolated manner prevents the use of similar sequences to achieve more accurate representations of the states. As for the computational complexity, these methods do not scale well with the dataset size $N$. Specifically, the number of likelihoods that need to be obtained is $N^2$ (or $PN$ using the KL method). In the case of PPKs, $N^2/2$ evaluations are required, since the kernel is symmetric.

## 3. State-space dynamics (SSD) distance

In this paper, we propose to take a different approach in order to overcome the need of fitting an HMM to each sequence. To this end, we propose to train a single, large HMM $\Theta$ of $K$ hidden states using all the sequences in the dataset. This will allow for a better estimation of the emission probabilities of the hidden states, compared to the case where an HMM is trained on each sequence. Then, we use the state-space of $\Theta$ as a common representation for the sequences. Each sequence $\mathbf{S}_n$ is linked to the common state-space through the transition matrix that it induces when is fed into the model. This matrix is denoted as $\tilde{\mathbf{A}}^n = \{\tilde{a}_{ij}^n\}_{i,j=1}^K$, where

$$\tilde{a}_{ij}^n = p(q_{t+1}^n = s_j|q_t^n = s_i, \mathbf{S}_n, \Theta). \tag{8}$$

In order to obtain each $\tilde{\mathbf{A}}^n$, we run the forward–backward algorithm for the sequence $\mathbf{S}_n$ under the parameters $\Theta$ (including the learned transition matrix $\mathbf{A} = \{a_{ij}\}$) and then obtain the sequence-specific transition probabilities by using Eq. (5):

$$\tilde{a}_{ij}^n \propto \sum_{t'=1}^T \alpha_i^n(t') a_{ij} p(\mathbf{x}_{t'+1}|q_{t'+1} = s_j) \beta_j^n(t'+1), \tag{9}$$

where $\alpha_i^n(t)$ and $\beta_j^n(t'+1)$ are the forward and backward variables for $\mathbf{S}_n$, respectively. This process can be seen as a projection of the dynamical characteristics of $\mathbf{S}_n$ onto the state-space defined by the common model $\Theta$. Therefore, the overall transition matrix $\mathbf{A}$ of the large model $\Theta$ acts as a common, data-dependent "prior" for the estimation of these individual transition matrices.

When using dynamical models for sequence clustering, the use of the proposed common set of emission distributions can be motivated as follows: if this kind of algorithms are used, the most likely situation is that there is a high degree of state-sharing between the models for the different classes. If that were not the case, the "static" probability densities of the different classes would hardly overlap and the dynamical information would not be required, so simpler models such as mixtures of Gaussians could be used.

This procedure is somewhat equivalent to obtaining individual HMMs with emission distributions that are shared or "clamped" amongst the different models. Clamping is a usual and useful tool when one wants to reduce the number of free parameters of a model in order to either obtain a better estimate or reduce the computational load. In our case, the effects of clamping the emission distributions are twofold: we get the usual benefit of better estimated parameters and, at the same time, it allows for simple distance measures between hidden Markov models using the transition distributions. This happens because the transition processes of the different models now share a common support, namely the fixed set of emission distributions. As previously mentioned, running the forward–backward algorithm implies a

time complexity of $O(K^2T)$ for a sequence of length $T$, which is the same complexity required for obtaining the likelihood of an HMM. Our proposal only requires $N$ of these calculations, instead of $N^2$ likelihood evaluations or $N^2/2$ PPKs as the methods mentioned in the previous section do. This makes the SSD algorithm a valuable method for working with large datasets.

At this point, we have each sequence $\mathbf{S}_n$ represented by its induced transition matrix $\tilde{\mathbf{A}}^n$. In order to define a meaningful distance measure between these matrices, we can think of each $\tilde{\mathbf{A}}^n = [\mathbf{a}_{n1}, \ldots, \mathbf{a}_{nK}]^T$ as a collection of $K$ discrete probability functions $\mathbf{a}_{n1}, \ldots, \mathbf{a}_{nK}$, one per row, corresponding with the transition probabilities from each state to every other state. In this manner, the problem of determining the affinity between sequences can finally be transformed into the well-studied problem of measuring similarity between distributions. In this work, we employ the Bhattacharyya affinity [17], defined as

$$D_B(p_1, p_2) = \sum_x \sqrt{p_1(x)p_2(x)}, \tag{10}$$

where $p_1$ and $p_2$ are discrete probability distributions. We consider the affinity between two transition matrices to be the mean affinity between their rows. The distance between two sequences $\mathbf{S}_i$ and $\mathbf{S}_j$ can then be written as

$$d_{ij}^{\text{BHAT}} = -\log \frac{1}{K} \sum_{k=1}^{K} D_B(p_{ik}, p_{jk}). \tag{11}$$

Other approaches could be used in order to define distances between the different transition matrices. For example, instead of using $\tilde{\mathbf{A}}^n$ directly, an idea similar to diffusion distances [18] could be applied by using different powers of the transition matrices $(\tilde{\mathbf{A}}^n)^t$, where $t$ is a time index. This is equivalent to iterating the random walk defined by the transition matrices for $t$ time steps. The $j$th row of such an iterated transition matrix encodes the probabilities of transitioning from state $j$ to each other state in $t$ time steps. However, this approach would introduce the extra parameter $t$, which must be set very carefully. For example, many transition matrices converge very quickly to the stationary distribution even for low $t$ (specially if the number of states is small). This could be a problem in cases where the stationary distributions for sequences in different clusters are the same. An example of such a scenario is presented in Section 4.

Moreover, the SSD distance measure is very flexible. Measuring distances between sequences is highly subjective and application dependent. For example, in a certain scenario we may not be interested in the rest time for each state, but only in the transitions (similar to Dynamic Time Warping [19]). To this end, a good alternative would be to obtain the transition matrices $\tilde{\mathbf{A}}^n$ for every sequence, but ignore the self-transitions in the distance measurement. That can be easily done by setting all the self-transitions to 0 and then re-normalizing the rows of the resulting transition matrices.

Once the distances between all the sequences are obtained, the actual clustering can be carried out using spectral clustering (or any other typical technique). We refer to this algorithm as state-space dynamics (SSD) clustering. It is summarized in Algorithm 1.

**Algorithm 1.** SSD distance for clustering sequential data.

*Inputs:*
Dataset $\mathcal{S} = \{\mathbf{S}_1, \ldots, \mathbf{S}_N\}$, $N$ sequences
$K$: Number of hidden states
*Algorithm:*
Step 1: Learning the global model (Baum–Welch)
$\quad \Theta = \arg\max_{\Theta'} P(\mathbf{S}_1, \ldots, \mathbf{S}_N | \Theta')$
Step 2: Estimating $\tilde{\mathbf{A}}^n = \{\tilde{a}_{ij}^n\}$ (Forward/Backward)
**for all $\mathbf{S}_n$ do**

$\quad \alpha_k(t) = P(\mathbf{S}_n(1), \ldots, \mathbf{S}_n(t), q_t = k | \Theta)$
$\quad \beta_k(t) = P(\mathbf{S}_n(t+1), \ldots, \mathbf{S}_n(T_n), q_t = k | \Theta)$
$\quad \tilde{a}_{ij}^n \propto \sum_{t=1}^{T_n} \alpha_i(t)a_{ij}p(\mathbf{S}_n(t+1)|q_{t+1}=i)\beta_j(t+1)$
**end for**
Step 3: Obtaining the distance matrix $\mathbf{D} = \{d_{ij}\}$
**for all $i,j$ do**

$\quad \mathbf{p}_{ik} \equiv k^{th}$ row of $\tilde{\mathbf{A}}^i$
$\quad d_{ij} = -\log(1/K) \sum_{k,k'=1}^{K} \sqrt{p_{ik}(k')p_{jk}(k')}$
**end for**
Step 4: Clustering using $\mathbf{D}$

It is worth noting that our proposal does not include any special initialization of the large model representing the dataset, such as imposing a block-diagonal structure on the transition matrix to encourage the clustering [6]. We do not aim to obtain a single generative model of the complete dataset, but an adequate common representation that allows for a subsequent successful non-parametric clustering.

An important free parameter of our method is $K$, the number of hidden states of the common model. It should be chosen accordingly to the richness and complexity of the dataset. In the worst case (that is to say, assuming that there is no state sharing amongst different groups), it should grow linearly with the number of groups. In the experiments included in this work, we have fixed this size a priori, but it could be estimated using well-known criteria such as BIC or AIC [20].

Nonetheless, we do not expect this parameter $K$ to be critical for the success of the method as long as it is within a sensible range, and we prove this point in the experiments by trying a wide range of hidden space cardinalities. Further discussion about this topic is provided in Section 4.3.

Recall that the forward–backward algorithm for HMMs is $O(K^2T)$, where $K$ is the number of states and $T$ the sequence length. This indicates that our proposal is specially suitable for datasets consisting of a large number of sequences coming from a small number of clusters, which is a usual case. In such a scenario, the number of hidden states required for a successful clustering is low, so the time penalty in the FW–BW algorithm will be more than compensated by the significant computational load reduction coming from the linear complexity in the number of sequences. If sequences coming from different clusters share some emission distributions, the improvements will be even more notorious, because the algorithm will exploit that sharing in a natural way.

As previously explained, dynamic model-based distance measures are specially useful when there is a significant overlap between the emission distributions of different classes. Our proposal is thus very suitable for working in this interesting scenario. If there were no such sharing, the model would have to be large enough to accommodate the individual models for all the classes. If the emission models for the different classes were distinct enough, a given sequence would only transition between the states representing its class and thus the distances between sequences belonging to different clusters would be infinite.

Finally, we would like to comment on the relationship between our proposal and that of [21]. There, the authors propose a Bayesian clustering method based on transition matrices of Markov chains. They assume that the sequences are discrete, so a Markov chain can be directly estimated via transition counts. Our proposal, on the other hand, uses Markov chains on the latent variables (states), what makes it far more general. Moreover, our focus is on defining a general model-based distance between sequences, so that the SSD distance can be directly coupled with a wide range of clustering algorithms depending on the task at hand.

# 4. Experimental results

In this section we present a thorough experimental comparison between SSD and state of the art algorithms using both synthetic and real-world data. Synthetic data include an ideal scenario where the sequences in the dataset are actually generated using HMMs, as well as a control chart clustering task. Real data experiments include different scenarios (character, gesture and speaker clustering) selected from the UCI-ML [22] and UCI-KDD [23] repositories. The implementation of the compared algorithms is provided in the author's website,[1] except for PPK, which is available at http://www1.cs.columbia.edu/~jebara/code.html.[2]

The compared methods for obtaining the distance matrix are: **SSD**, state-space dynamics clustering with Bhattacharyya distance; **PPK**, probability product kernels [11]; **KL**, KL-divergence based distance [7]; **BP**, BP metric [8]; **YY**, Yin–Yang distance [10] and **SYM**, symmetrized distance [6].

We denote the number of hidden states of the global model used by SSD as $K$, and the number of states per model of the methods that rely on training a HMM on each single sequence as $K_m$.

Once a distance matrix is available, we perform the actual clustering using the spectral algorithm described in [2]. The different distance matrices are turned into similarity matrices by means of a Gaussian kernel whose width is automatically selected in each case attending to the eigengap. Though more elaborated methods such as [25] can be used to select the kernel width, in our experiments it is automatically selected in each case attending to the eigengap since the experimental results are good enough. We assume that the number of clusters is known a priori. If this is not the case, automatic determination of the number of clusters can be carried out by methods such as those in [24,25]. The PPK method directly returns a similarity matrix, that is first converted into a distance matrix by taking the negative logarithm of each element. Then, it is fed into the clustering algorithm with automatic kernel width selection. The final $k$-means step of the spectral clustering algorithm is run 10 times, choosing as the final partition the one with the minimum intra-cluster distortion. The free parameter $T_{PPK}$ of the PPK method is fixed to 10 following [11].

The results shown in the sequel are averaged over a number of iterations in order to account for the variability coming from the EM-based training of the HMM. Performance is measured in the form of clustering accuracy, understood as the percentage of correctly classified samples under an optimal permutation of the cluster labels, or its reciprocal, the clustering error.

## 4.1. Synthetic data

In this subsection we test the algorithms using two kinds of synthetically generated data: a mixture-of-HMMs (MoHMM) scenario as in [6,7], and a UCI-ML dataset representing control charts.

### 4.1.1. Mixture of HMMs

Each sequence in this dataset is generated by a mixture of two equiprobable HMMs $\theta_1$ and $\theta_2$. Each of these models has two hidden states, with an uniform initial distribution, and their corresponding transition matrices are

$$\mathbf{A}_1 = \begin{pmatrix} 0.6 & 0.4 \\ 0.4 & 0.6 \end{pmatrix}, \quad \mathbf{A}_2 = \begin{pmatrix} 0.4 & 0.6 \\ 0.6 & 0.4 \end{pmatrix}.$$

Emission probabilities are the same in both models, specifically $N(0,1)$ in the first state and $N(3,1)$ in the second. This is a deceptively simple scenario. Since both the emission probabilities and the equilibrium distributions are identical for both models, the

only way to differentiate sequences generated by each of them is to attend to their dynamical characteristics. These, in turn, are very similar, making this a hard clustering task. The length of each individual sequence is uniformly distributed in the range $[0.6\mu_L, 1.4\mu_L]$, where $\mu_L$ is the mean length. Since the number of sequences coming from each model is the same, a trivial (assign all elements to the same cluster) clustering scheme would achieve a 50% baseline clustering error.

Fig. 1 shows the clustering error achieved by the compared algorithms in a dataset of $N=100$ sequences, averaged over 50 runs, for mean sequence lengths ranging from 25 to 250. All the algorithms based on individual models use the correct model structure ($K_m=2$ hidden states per class) to fit each sequence. For SSD, this implies using four hidden states for the common model ($K=4$). As expected, when the sequences follow an HMM generative model and the representative model structure is chosen accordingly, SSD achieves impressive performance improvements for short sequence lengths. In contrast, algorithms that rely on training an HMM for each sequence suffer from poor model estimation when the mean sequence length is very low ($\leq 100$), which in turn produces bad clustering results. Our proposal overcomes this difficulty by using information from all the sequences in order to generate the common representative HMM. Consequently, the emission probabilities are estimated much more accurately and the distances obtained are more meaningful, leading to a correct clustering. Nonetheless, when the sequences are long ($\geq 200$) very accurate models can be obtained from each single sequence and the different methods tend to converge in performance.

### 4.1.2. Synthetic control chart

This dataset contains unidimensional time series representing six different classes of control charts: normal, cyclic, increasing trend, decreasing trend, upward shift and downward shift. There are 100 instances of each class, with a fixed length of 60 samples per instance. A sample of each class is plotted in Fig. 2.

We carry out a multi-class clustering task on this dataset, partitioning the corpus into six groups which we expect to correspond with the different classes of control charts. As explained in Section 3, the size of the state-space for the HMM in SSD clustering should be chosen accordingly to the number of classes, so we employ a number of hidden states in the range 12–40. It also allows us to show that our proposal is robust enough to produce good performance in such an extensive range. Results, averaged over 10 runs, are shown in Table 1.
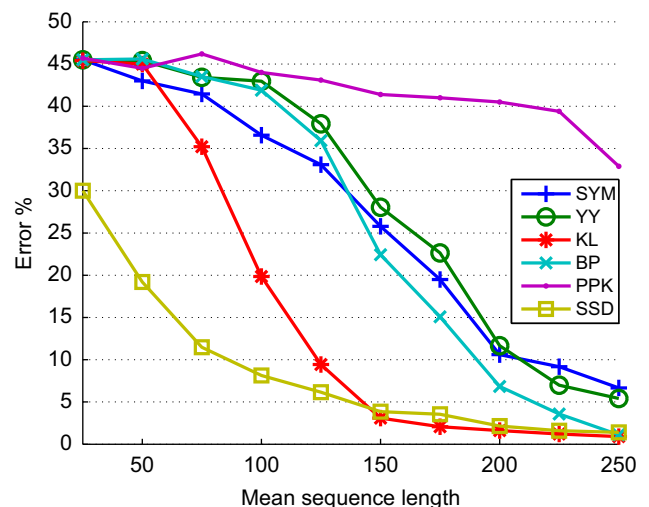


**Fig. 1.** Clustering error for the MoHMM case.

In general, all methods obtain good results, specially taking into account the multi-class nature of the dataset (the baseline trivial clustering assigning all sequences to the same cluster would attain a 16.66% accuracy). However, SSD clearly outperforms all the other compared algorithms. It is specially remarkable the very high performance achieved in the "large-dataset" of 100 sequences per class. Such good results are obtained because, in contrast to previous proposals, the modeling employed by SSD distance improves as the dataset size increases.
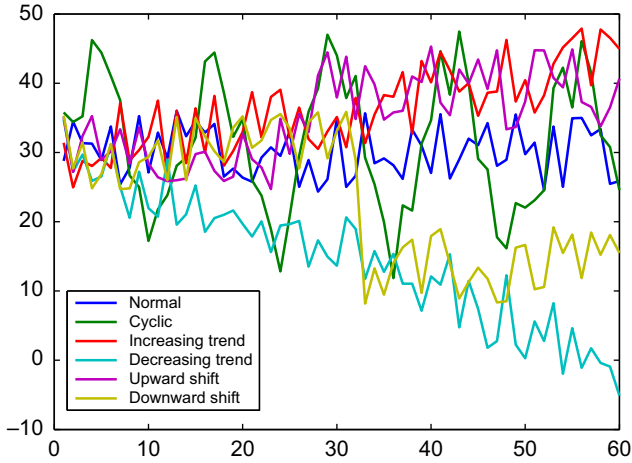


**Fig. 2.** Some samples from the synthetic control chart dataset.

It is worth noting how SSD provides better results than any other compared method even for the very modest number of 12 hidden states. This implies that there is a great amount of state-sharing between the different classes. Also of interest is the fact that the performance gets even better as that number of states is increased. We tried to push the limits of the method by using the largest number of hidden states that could be handled in a reasonable time (40 states), and even in that case we obtained very good results. This supports the idea that the size of the common model space is not a critical parameter of the algorithm. The confusion matrix when $N=30$ and $K=20$ (averaged over the 10 runs) is shown in Fig. 3 in the form of a Hinton diagram.

### 4.2. Real-world data clustering experiments

We use the following datasets from the UCI ML and KDD archives:

*Character Trajectories*: This dataset consists of trajectories captured by a digitizing tablet when writing 20 different characters. Each sample is a three-dimensional vector containing the $x$ and $y$ coordinates as well as the pen tip force. The sequences are already differentiated and smoothed using a Gaussian kernel. We use 25 sequences per character, and carry out two-class clustering between all the possible combinations, giving a total of 190 experiments. The average length of the sequences in this dataset is around 170 samples. The baseline error is 50%.

*AUSLAN*: The Australian Sign Language dataset is comprised 22-dimensional time series representing different sign-language

**Table 1**
Mean accuracy (standard deviation in brackets) in the control chart dataset. Using (a) 30 and (b) 100 sequences per class.

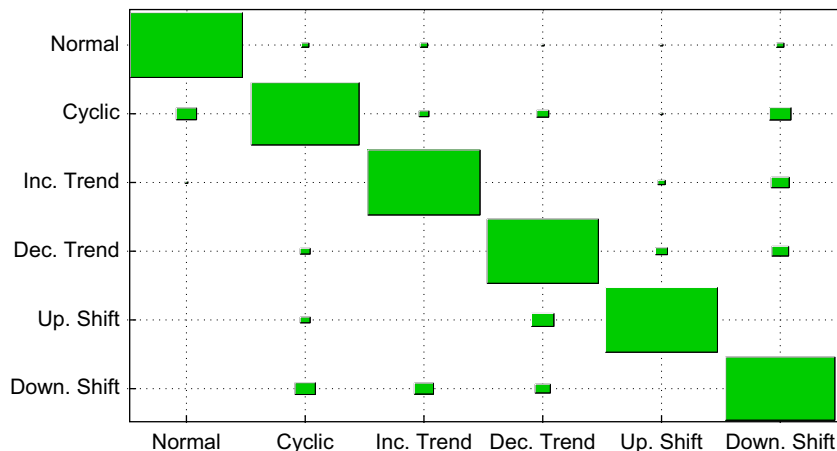| $K_m$ | SYM | YY | KL | BP | PPK | $K$ | SSD |
|---|---|---|---|---|---|---|---|
| (a) 30 sequences per class | | | | | | | |
| $K_m=2$ | 76.11% ($\pm$0.1) | 74.11% ($\pm$5.3) | 74.67% ($\pm$5.3) | 78.33% ($\pm$0.2) | 46.89% ($\pm$3.3) | $K=12$ | 81.61% ($\pm$6.0) |
| $K_m=3$ | 74.00% ($\pm$0.6) | 74.67% ($\pm$0.5) | 78.67% ($\pm$3.9) | 75.89% ($\pm$3.6) | 55.22% ($\pm$0.6) | $K=16$ | 86.28% ($\pm$6.3) |
| $K_m=4$ | 74.78% ($\pm$0.6) | 76.44% ($\pm$1.6) | 79.33% ($\pm$3.1) | 76.00% ($\pm$4.4) | 51.22% ($\pm$4.8) | $K=20$ | 87.33% ($\pm$4.1) |
| $K_m=5$ | 77.00% ($\pm$1.8) | 79.11% ($\pm$3.2) | 77.44% ($\pm$4.9) | 79.78% ($\pm$1.6) | 41.78% ($\pm$3.1) | $K=28$ | **88.19%** ($\pm$4.7) |
| $K_m=6$ | 74.67% ($\pm$0.3) | 76.89% ($\pm$3.1) | 76.22% ($\pm$2.9) | 74.44% ($\pm$2.0) | 40.11% ($\pm$3.5) | $K=40$ | 85.83% ($\pm$2.7) |
| (b) 100 sequences per class | | | | | | | |
| $K_m=2$ | 67.66% ($\pm$4.35) | 73.30% ($\pm$7.1) | 70.19% ($\pm$0.08) | 72.91% ($\pm$4.9) | 45.79% ($\pm$1.8) | $K=12$ | 91.57% ($\pm$1.5) |
| $K_m=3$ | 79.23% ($\pm$0.26) | 85.53% ($\pm$3.2) | 70.2% ($\pm$0.07) | 87.33% ($\pm$0.6) | 51.23% ($\pm$1.9) | $K=16$ | 92.71% ($\pm$1.8) |
| $K_m=4$ | 79.28% ($\pm$0.7) | 87.15% ($\pm$3.6) | 72.23% ($\pm$3.0) | 84.92% ($\pm$4.3) | 53.15% ($\pm$3.7) | $K=20$ | 93.23% ($\pm$1.4) |
| $K_m=5$ | 81.12% ($\pm$3.4) | 84.66% ($\pm$4.9) | 70.90% ($\pm$0.5) | 82.35% ($\pm$2.5) | 40.75% ($\pm$5.8) | $K=28$ | **94.07%** ($\pm$1.1) |
| $K_m=6$ | 78.63% ($\pm$5.5) | 82.78% ($\pm$4.1) | 73.95% ($\pm$4.5) | 80.73% ($\pm$5.2) | 39.15% ($\pm$1.7) | $K=40$ | 93.60% ($\pm$1.0) |



**Fig. 3.** Confusion matrix for SSD clustering with 30 sequences per class and 20 hidden states.

gestures. The gestures belong to a single signer, and were collected in different sessions over a period of nine weeks. There are 27 instances per gesture, with an average length of 57 samples. Following [11], we perform 2-class clustering tasks using semantically related concepts. These concepts are assumed to be represented by similar gestures and thus provide a difficult scenario. Baseline error is once again 50%.

*Japanese Vowels*: To construct this dataset, nine male speakers uttered two Japanese vowels consecutively. The actual data is comprised the 12-dimensional time-series of LPC cepstrum coefficients for each utterance, captured at a sampling rate of 10 kHz using a sliding window of 25.6 ms with a 6.4 ms shift. The number of samples per sequence varies in the range 7–29 and there are 30 sequences per user. We use this dataset for two different tasks: speaker clustering and speaker segmentation. Since it is a 9-class clustering task with the same number of sequences per class, the trivial clustering error baseline is 11.11%.

Table 2 shows the numerical results, averaged over 10 runs. In the Character Trajectories dataset, the individual sequences are fairly long and the classes are mostly well separated, so this is an easy task. Single sequences are informative enough to produce good representative models and, consequently, most methods achieve very low error rates. Nonetheless, using the SSD distance outperforms the competitors.

For the AUSLAN dataset, following [11], we used HMMs with $K_m=2$ hidden states for the methods that train a single model per sequence. The sequences were fed directly to the different algorithms without any preprocessing. We reproduce the results for the PPK method from [11]. All methods perform well over the baseline, except for the symmetric distance. The common model for SSD employs $K=4$ hidden states ($2K_m$), since the 2-way clustering tasks are fairly simple in this case. It is worth noting that the bad performance in the 'Yes' vs 'No' case is due to a dataset artifact that causes the algorithms to try to cluster the sequences attending to the recording session instead of to the actual sign they represent. Our proposal produces great results in this dataset, surpassing the

rest of the methods in every pairing except for the pathological 'Yes' vs 'No' case.

Finally, we carry out a 9-class speaker clustering task using the full Japanese Vowels dataset. There are 30 sequences per class. The task is not too complex, as all the methods perform well over the baseline. The large number of classes and their variability demands a large number of hidden states in the common HMM of SSD. This, in turn, means a time penalty as the HMM training time is quadratic in the number of hidden states. Nonetheless, the performance obtained in this dataset by our proposal is very competitive in terms of clustering accuracy, only being surpassed by the KL method. It is also remarkable how the SSD-based clustering exhibits a very stable performance in a huge range of state-space cardinalities, what once again coincides with our intuition that an accurate determination of that parameter is not crucial to the algorithm.

### 4.2.1. Speaker segmentation

In order to briefly show another application of sequence clustering, we have reproduced the Japanese Vowels dataset segmentation experiment from [26] using our proposed SSD distance. This scenario is constructed by concatenating all the individual sequences in the dataset to form a long sequence which we would like to divide into nine segments (one per user). Only two methods, KL and SSD have been tested in this task because they were the best-performing distances for clustering this dataset. In order to carry out the sequence-clustering-based segmentation, we first extract subsequences using a non-overlapping window. The length of these windows range from 10 to 20 samples. When using the KL distance, each subsequence is modelled using a 2-state HMM, which is the optimal value for the corresponding clustering task. A distance matrix is then obtained from these subsequences using the KL or SSD distances, and then spectral segmentation (SS) is carried out instead of spectral clustering [26]. This amounts to performing dynamic programming (DP) on the eigenvectors

**Table 2**
Performance on the Character Trajectories (top), AUSLAN (middle) and Japanese Vowels (bottom, 9-class clustering task) datasets. The standard deviation of the results for the AUSLAN dataset is 0 in every case except for 'SPEND' vs 'COST' using YY distance, with a value of 0.8. The number of hidden states is $K=4$ for SSD and $K_m=2$ for the rest of methods (best case).

| # hidden stat. | SYM | YY | KL | BP | PPK | # hidden stat. | SSD |
|---|---|---|---|---|---|---|---|
| $K_m=2$ | 96.10% | 97.02% | 96.42% | 96.57% | 76.72% | $K=14$ | 97.17% |
| | ($\pm 0.4$) | ($\pm 0.2$) | ($\pm 0.1$) | ($\pm 0.2$) | ($\pm 0.8$) | | ($\pm 0.3$) |
| $K_m=3$ | 96.30% | 96.90% | 96.45% | 95.23% | 67.66% | $K=16$ | 97.58% |
| | ($\pm 0.1$) | ($\pm 0.2$) | ($\pm 0.0$) | ($\pm 0.3$) | ($\pm 0.8$) | | ($\pm 0.2$) |
| $K_m=4$ | 95.31% | 95.53% | 95.69% | 83.92% | 62.25% | $K=20$ | 98.23% |
| | ($\pm 0.2$) | ($\pm 0.0$) | ($\pm 0.1$) | ($\pm 0.8$) | ($\pm 0.2$) | | ($\pm 0.2$) |
| $K_m=5$ | 96.28% | 96.40% | 96.58% | 84.70% | 61.39% | $K=22$ | **98.35%** |
| | ($\pm 0.3$) | ($\pm 0.1$) | ($\pm 0.1$) | ($\pm 0.1$) | ($\pm 0.2$) | | ($\pm 0.2$) |
| SIGNS | SYM (%) | YY (%) | KL (%) | BP (%) | PPK (%) | | SSD (%) |
| 'HOT' vs 'COLD' | **100** | **100** | **100** | **100** | **100** | | **100** |
| 'EAT' vs 'DRINK' | 51.85 | 92.59 | 92.59 | 92.59 | 93 | | **95.37** |
| 'HAPPY' vs 'SAD' | 59.26 | 98.15 | **100** | 98.15 | 87 | | **100** |
| 'SPEND' vs 'COST' | 54.07 | 99.63 | **100** | **100** | 80 | | **100** |
| 'YES' vs 'NO' | **60.36** | 55.56 | 55.56 | 55.56 | 59 | | 56.66 |
| # hidden stat. | SYM | YY | KL | BP | PPK | # hidden stat. | SSD |
| $K_m=2$ | 66.67% | 85.11% | **90.15**% | 85.30% | 75.48% | $K=20$ | 82.30% |
| | ($\pm 2.8$) | ($\pm 2.1$) | ($\pm 1.0$) | ($\pm 1.0$) | ($\pm 3.7$) | | ($\pm 3.7$) |
| $K_m=3$ | 67.18% | 79.01% | 81.14% | 75.44% | 82.59% | $K=30$ | 85.48% |
| | ($\pm 3.0$) | ($\pm 4.1$) | ($\pm 4.0$) | ($\pm 3.7$) | ($\pm 5.1$) | | ($\pm 4.5$) |
| $K_m=4$ | 67.96% | 83.41% | 85.55% | 78.55% | 79.78% | $K=40$ | 87.93% |
| | ($\pm 3.4$) | ($\pm 6.2$) | ($\pm 5.4$) | ($\pm 4.8$) | ($\pm 3.7$) | | ($\pm 4.4$) |
| $K_m=5$ | 70.44% | 82.81% | 82.30% | 79.77% | 78.15% | $K=50$ | 86.37% |
| | ($\pm 3.6$) | ($\pm 5.2$) | ($\pm 5.5$) | ($\pm 4.4$) | ($\pm 3.9$) | | ($\pm 6.9$) |

**Table 3**
Segmentation error (mean and standard deviation) in the Japanese Vowels dataset (nine sources and segments) using KL and SSD distances.

| Window length | KL-SS | SSD-SS | | | |
|---|---|---|---|---|---|
| | $K_m=2$ | $K=18$ | $K=24$ | $K=32$ | $K=40$ |
| $W=10$ | 1.0% ($\pm 0.22$) | 1.61% ($\pm 0$) | 1.0% ($\pm 0.29$) | 0.82% ($\pm 0.35$) | 0.72% ($\pm 0.33$) |
| $W=15$ | 3.5% ($\pm 0$) | 2.39% ($\pm 0$) | 1.12% ($\pm 0.39$) | 0.95% ($\pm 0.29$) | 0.98% ($\pm 0.27$) |
| $W=20$ | 1.4% ($\pm 0$) | 2.66% ($\pm 0$) | 1.26% ($\pm 0.44$) | 1.03% ($\pm 0.53$) | 0.75% ($\pm 0.45$) |

resulting from the decomposition of the normalized Laplacian matrix used for spectral clustering. In the ideal case, those eigenvectors are expected to be approximately piece-wise constant on each cluster, thus being a very adequate representation to run DP on in order to obtain a segmentation of the subsequences.

Table 3 shows the performance obtained in this task. Performance for the KL method is taken directly from [26]. Segmentation error is measured as the number of incorrectly "classified" segments in the sequence, which can be seen as the fraction of the time that the segmentation algorithm is giving wrong results. As the results confirm, using the SSD distance is very adequate in this scenario because the number of subsequences is quite large and, at the same time, all of them are very short. This way, we exploit both the reduced time complexity in the number of sequences and the better estimation of the emission distributions.

### 4.3. On the number of hidden states for SSD distance

A general conclusion that can be drawn from the experimental results is that the proposed distance measure generally increases its performance as the common state-space grows larger. This way, the size of the initial HMM can be chosen according to the available computational power. A performance loss is predictable at some point, although we have not witnessed it in our experiments, since there the limiting factor have always been the computational load. A possible way to test if the number of hidden states is getting too large for the problem at hand is to look at the mean occupancy (or stationary probability) of those states. If there are some of them which are occupied only for a very small fraction of the time, this can be seen as a sign that the common model is starting to overfit to particular characteristics of some individual sequences. In this case, the state space is too large and should be reduced.

## 5. Conclusions

In this paper we have presented a new distance for model-based sequence clustering using state-space models. We learn a single model representing the whole dataset and then obtain distances between sequences attending to their dynamics in the common state-space that this model provides. It has been empirically shown that the proposed approach outperforms the previous semi-para-metric methods, specially when the mean sequence length is short. Furthermore, the proposed method scales much better with the dataset size (linearly vs quadratically). As drawback of this method it should be mentioned that, as the number of classes grow, the common model may need a large number of hidden states to correctly represent the dataset (although the method is empirically shown not to be too sensitive to the accurate determination of the model size). In the case of hidden Markov models, the time complexity of the training procedure is quadratic in this number of states, so the total running time can be high in these cases. Consequently, we find our proposal specially appealing for scenarios with a large number of sequences coming from a few different classes, which is a very usual case.

Promising lines for future work include the application of this methodology to other state-space models, both discrete and continuous, and to semi-supervised scenarios. We are also investigating alternative definitions of distance measures between transition matrices in order to take into account the potential redundancy of the state-space.

## Acknowledgments

## References

[1] R. Xu, D. Wunsch II, Survey of clustering algorithms, IEEE Trans. Neural Networks 16 (3) (2005) 645–678.

[2] A. Ng, M. Jordan, Y. Weiss, On spectral clustering: analysis and an algorithm, in: Advances in Neural Information Processing Systems, 2002.

[3] U. von Luxburg, A tutorial on spectral clustering, Statistics and Computing 17(4) (2007).

[4] T.W. Liao, Clustering of time series data—a survey, Pattern Recognition 38 (11) (2005) 1857–1874. doi:10.1016/j.patcog.2005.01.025 URL: ⟨http://www.sciencedirect.com/science/article/B6V14-4G7G4BP-2/2/8e96ce614dd47a4a84c6712bd0c43022⟩.

[5] L. Rabiner, A tutorial on hidden Markov models and selected applications in speech recognition, Proc. IEEE 77 (2) (1989) 257–286.

[6] P. Smyth, Clustering sequences with hidden Markov models, in: Advances in Neural Information Processing Systems, vol. 9, 1997, pp. 648–654.

[7] D. García-García, E. Parrado-Hernández, F. Díaz-de-María, A new distance measure for model-based sequence clustering, IEEE Trans. Pattern Anal. Mach. Intell. 31 (7) (2009) 1325–1331. doi:10.1109/TPAMI.2008.268.

[8] A. Panuccio, M. Bicego, V. Murino, A hidden Markov model-based approach to sequential data clustering, in: Proceedings of the Joint IAPR International Workshop on Structural, Syntactic and Statistical Pattern Recognition, 2002, pp. 734–742.

[9] F. Porikli, Clustering variable length sequences by eigenvector decomposition using HMM, in: Proceedings of the International Workshop on Structural and Syntactic Pattern Recognition, Lisbon, Portugal, 2004, pp. 352–360.

[10] J. Yin, Q. Yang, Integrating hidden Markov models and spectral analysis for sensory time series clustering, in: Fifth IEEE International Conference on Data Mining, 2005.

[11] T. Jebara, Y. Song, K. Thadani, Spectral clustering and embedding with hidden Markov models, in: Proceedings of the 18th European Conference on Machine Learning (ECML), Warsaw, Poland, 2007.

[12] J. Alon, S. Sclaroff, G. Kollios, V. Pavlovic, Discovering clusters in motion time-series data, in: Proceedings of the 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'03), 2003.

[13] T. Oates, L. Firoiu, P. Cohen, Using dynamic time warping to bootstrap HMM-based clustering of time series, in: Sequence Learning—Paradigms, Algorithms, and Applications, Springer-Verlag, London, UK, 2001, pp. 35–52.

[14] J. Shi, J. Malik, Normalized cuts and image segmentation, IEEE Trans. Pattern Anal. Mach. Intell. 22 (8) (2000) 888–905.

[15] M. Belkin, P. Niyogi, Laplacian eigenmaps and spectral techniques for embedding and clustering, in: Advances in Neural Information Processing Systems 14, MIT Press, 2001, pp. 585–591.

[16] G. Golub, C. Van Loan, Matrix Computations, John Hopkins University Press, 1989.

[17] A. Bhattacharyya, On a measure of divergence between two statistical populations defined by their probability distributions, Bull. Calcutta Math Soc., 1943.

[18] A. Szlam, R. Coifman, M. Maggioni, A general framework for adaptive regularization based on diffusion processes, J. Mach. Learn. Res. 9 (9) (2008) 1711–1739.

[19] H. Sakoe, S. Chiba, Dynamic programming algorithm optimization for spoken word recognition, IEEE Trans. Acoust. Speech Signal Process. 26 (1) (1978) 43–49.
[20] C.M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.
[21] M. Ramoni, P. Sebastiani, P. Cohen, Bayesian clustering by dynamics, Mach. Learn. 47 (1) (2002) 91–121 doi: http://dx.doi.org/10.1023/A:1013635829250.
[22] C. Blake, C. Merz, UCI Repository of machine learning databases, University of California, Irvine, Department of Information and Computer Sciences ⟨http://www.ics.uci.edu/~mlearn/MLRepository.html⟩.
[23] S. Hettich, S. Bay, The UCI KDD Archive, University of California, Irvine, Department of Information and Computer Science ⟨http://kdd.ics.uci.edu⟩.
[24] G. Sanguinetti, J. Laidler, N.D. Lawrence, Automatic determination of the number of clusters using spectral algorithms, in: IEEE Workshop on Machine Learning for Signal Processing, 2005, pp. 55–60. doi: 10.1109/MLSP.2005.1532874. URL ⟨http://dx.doi.org/10.1109/MLSP.2005.1532874⟩.
[25] L. Zelnik-Manor, P. Perona, Self-tuning spectral clustering, Advances in Neural Information Processing Systems, 2004, pp. 1601–1608.
[26] D. García-García, E. Parrado-Hernández, F.D. de María, Sequence segmentation via clustering of subsequences, in: International Conference on Machine Learning and Applications (ICMLA), 2009.

**Darío García-García** earned his Telecommunications Engineering degree in 2006 from University de Cantabria, Spain. He is currently working towards the Ph.D. degree in the Multimedia Processing Group of Universidad Carlos III de Madrid. His research interests include statistical learning theory, unsupervised learning, feature selection and working with sequential data.

**Emilio Parrado-Hernández** received his Telecommunications Engineering degree from the University of Valladolid, Spain in 1999 and his Doctorate from the University Carlos III of Madrid, Spain in 2003. He is currently a lecturer with the Department of Signal Processing and Communications at the University Carlos III of Madrid. His research is focused on Machine Learning, specifically on kernel methods, such as Support Vector Machines and Spectral Clustering, as well as on the application of this techniques to signal and data processing problems.

**Fernando Díaz-de-María** received the Telecommunication Engineering degree in 1991 and the Ph.D. degree in 1996 from the Universidad Politécnica de Madrid, Spain. From October 1991 to October 1995, he worked as a Lecturer at Universidad de Cantabria, Santander, Spain. From November 1995 to September 1996, he worked as an Assistant Professor at Universidad Politécnica de Madrid. He reached the Ph.D. Degree in July 1996. His Ph.D. Thesis focused on speech coding. In particular, on the use of Artificial Neural Network-based nonlinear prediction in CELP ("Code-Excited Nonlinear Predictive")-type coders. From October 1996, he is an Associate Professor at the Department of Signal Processing & Communications, Universidad Carlos III de Madrid, Madrid, Spain. From October 97 till nowadays, he has held several offices in both his Department and his University. His primary research interests include robust speech recognition, video coding and multimedia indexing. He has led numerous Projects and Contracts in the mentioned fields. He is co-author of several papers in prestigious international journals, two chapters in international books and quite a few papers in revised national and international conferences.