# Spatio-Temporal Convolution Kernels

Konstantin Knauf  $\,\cdot\,$  Daniel Memmert  $\,\cdot\,$  Ulf Brefeld

Received: date / Accepted: date

**Abstract** Trajectory data of simultaneously moving objects is being recorded in many different domains and applications. However, existing techniques that utilise such data often fail to capture characteristic traits or lack theoretical guarantees. We propose a novel class of spatio-temporal convolution kernels (STCKs) to capture similarities in multi-object scenarios. The abstract kernel is a composition of a temporal and a spatial kernel and its actual instantiations depend on the application at hand. Empirically, we compare our kernels and efficient approximations thereof to baseline techniques for clustering tasks using artificial and real world data from team sports.

Keywords Convolution Kernel  $\cdot$  Spatio-Temporal  $\cdot$  Trajectory  $\cdot$  Soccer

# **1** Introduction

Trajectory data of simultaneously moving objects is the key to analyse animal migration [38], transportation [3,14], tactics in team sports [20,33,43,61], players and avatars in (serious) computer games [31,48], customer behaviour [35] as well as spread patterns of fires [57]. A characteristic trait of many such applications is that trajectories of several objects are more informative than the trajectory of a single object. For instance, a single trajectory of a bird is not indicative for bird migration as individuals may join or leave the flock [38]

Konstantin Knauf

TU Darmstadt, Knowledge Mining & Assessment Group E-mail: knauf.konstantin@gmail.com

Daniel Memmert German Sport University Cologne, Institute of Cognitive and Team Sport Research E-mail: memmert@dshs-koeln.de

Ulf Brefeld

TU Darmstadt and DIPF Frankfurt/Main, Knowledge Mining & Assessment Group E-mail: brefeld@cs.tu-darmstadt.de

and a single trajectory of a soccer player does not reveal insights on the actual situation on the pitch [15,61].

Therefore, trajectories of multiple objects need to be processed together. Although this insight sounds trivial, processing multiple trajectories simultaneously challenges the standard model of computation as trajectories may interdepend in time and space in multiple ways. To exploit these dependencies, it is necessary to establish a notion of similarity for spatio-temporal paths of multiple objects to identify frequent patterns. By definition, frequent patterns are formed by an *a priori* unknown subset of objects at unknown locations in time and space. Analysing multi-trajectory data is therefore inherently a combinatorial problem that involves processing data at large scales.

A second problem arises from existing methods for analysing spatio-temporal data. Traditional approaches often cannot deal with continuous spatial domains but rely on an appropriate discretisation of the data at-hand [30, 44,45]. However, finding an *a priori* optimal discretisation is often difficult in many domains where only the final result allows conclusions on whether an initial set of atomic events is plausible or not [30]. Furthermore, many approaches cannot deal with permutations of the objects and differences in speed, while still being sensitive to differences in the direction of the motion [20,28,61].

We devise a novel class of convolution kernels for multi-trajectory data. It is specially tailored to multi-object scenarios, i.e. trajectories of multiple simultaneously moving objects. The kernel properties as well as the modular nature of the proposed class of kernels renders it highly adaptive to different applications. Since it is a kernel, it can also naturally be deployed with any kernel machine. The three characteristics, *multi-object scenario, modularity* and *kernel property*, distinguish our approach from existing methods. Due to its distinct characteristics, our approach is more suitable for a large variety of applications, it is flexible with respect to the notion of similarity, and it is theoretically better grounded than most of the existing methods. Since the complexity of a kernel evaluation is quadratic in the number and the lengths of the involved objects, we also propose an efficient percental approximation. Empirically, the method is evaluated on artificial datasets and real-world tracking data from ten Bundesliga soccer matches. We generally observe that our convolution kernels lead to better clusterings compared to baseline methods.

The remainder of this article is structured as follows. Section 2 reviews existing work. Section 3 introduces our spatio-temporal convolution kernel methods and Section 4 reports on empirical results. Section 5 concludes.

## 2 Related Work

#### 2.1 Trajectory Clustering

Trajectory clustering, or clustering of spatio-temporal data respectively, has been an active field of research in the past years. Existing approaches mainly focus on the application of video surveillance with the goal to detect anomalies in the data stream [2,13,21,26,28,53]. Other applications include automatic sports analysis, weather evolution modelling, animal migration and traffic analysis. Existing approaches rely mostly on processing of single trajectories. Recent contributions in this area can be roughly grouped into *similarity-based approaches* [9,27,13,20,21,28,50] and *motion-based approaches* [2,59,26,39, 42,53].

Similarity-based approaches define pairwise similarities between trajectories which are then processed by some clustering algorithm. Junejo et al. [28] represent trajectories as a set of two-dimensional coordinates together with the Hausdorff distance. Subsequently, graph-cuts are deployed to recursively partition the trajectories. Hausdorff distances are also used to cluster trajectories by Chen at al. [27] where not only the position but also the direction of the trajectories is taken into account by using 4-tuples (x, y, dx, dy) instead of coordinates only. Fu et al. [13] first resample trajectories to obtain constant between-point distances. Then the corresponding points of two trajectories are compared using an RBF kernel where the longer trajectory is cut to the length of the shorter one. Spectral clustering is then used together with a symmetric normalised Laplacian.

Buzan et al. [9] extend the longest common subsequence algorithm to threedimensional coordinates and use a modified version of agglomerative hierarchical clustering. Hirano et al. [20] deploy multi-scale matchings to compare trajectories. The basic idea is to generate trajectories at different scales as convolutions of the trajectory and Gaussian kernels with different standard deviations. Their similarity measure is then based on the hierarchical structure of the trajectory segments at different scales. Subsequently, a rough clustering is employed. Piciarelli et al. [50] define a trajectory-to-cluster similarity by the average Euclidean distance of trajectory coordinates to the nearest cluster coordinate where offsets in time induce negative weights.

Our approach also belongs to these similarity-based methods. In general, similarity-based approaches suffer from two major drawbacks. First, their computational complexity is at least in quadratic in the number of trajectories. Second, they rely on clustering full trajectories and are hence sensitive to tracking errors and sub-trajectories. While the first drawback is inherent to all similarity-based methods, our distribution-based approach and gradual weighting mitigates the effects of noise and tracking errors and is able to identify partial matchings between trajectories.

In contrast to similarity-based approaches, *motion-based approaches* focus on local movements of objects to derive models for the overall (group) motion in a scene. Grimson et al. [59] and Jeong et al. [26] represent a trajectory by bags of positions as well as directions based on the bag-of-words representation of documents in natural language processing. To this end, the spatial domain is discretised and the number of occurrences of each position in a trajectory is counted. Grimson et al. also take into account temporal information by counting the occurrences of each (discretised) direction in a trajectory. The topic model Dual-HDP [59] is used to find semantic regions, which are combined to form the different trajectories. Jeong et al. use latent dirichlet allocation [8] to obtain semantic regions. To incorporate temporal information, a hidden Markov model is trained for each topic based on the sequences which are *close* to the topic. Saleemi et al. [53] propose kernel density estimation to learn a five dimensional distribution of transitions from  $(x_1, y_1)$  to  $(x_2, y_2)$  in time t. Markov chain monte carlo [1] is then deployed to sample the most likely paths given the learned transition probabilities.

Basharat et al. [2] also learn a model for transition probabilities. Instead of kernel density estimation, a Gaussian mixture model is fitted to the observed transitions. Lin et al. [42] exploit the Lie algebraic structure of affine transformations to learn a flow model consisting of overlapping two-dimensional Gaussian distributions, each of which corresponds to an affine transform dominant in this spatial area. The approach is applied to pedestrians in a train station and optical flows obtained from satellite images. Li et al. [39] also use a similar Lie algebraic representation called *spatial hybrid driving force model*, which, opposed to [42], evolves over time. This model is used to solve the so-called group motion segmentation problem, i.e. to answer the question of which objects take part in an organised group motion and which do not.

Motion-based approaches also inhere some limitations. First, they often neglect temporal information at least of second order (curvature). Second, they do not provide a mapping of the input trajectories to groups of similar trajectories but rather describe the combined motion of all objects in all trajectories over time. Our approach differs methodologically from the summarised techniques in several ways: First, it provides a general framework that covers many applications and properties as opposed to being a very specific similarity measure tailored to a single application domain. Second, our approach is specialised on multiple simultaneously moving objects instead of focussing on only trajectories of single objects. Third, being a kernel the similarity measure is straightforwardly applicable to a broad range of algorithms and is theoretically well grounded in contrast to heuristic approaches.

## 2.2 Sports Analytics

Current approaches in the area of sports game trajectory analysis either aim to define objective performance measures for players [29], classify [5,19,24,15, 49,55] or cluster [20,61] plays/trajectories, or learn a motion model for team behaviour [6,11,34,40,39,43,63].

In [29], Kang et al. define performance metrics for soccer players based on the definition of owned and competitive regions of the field, which are derived from ball and player trajectories. Siddiquie et al. [55] represent videos of football plays by a bag-of-features from histograms of optical flows as well as histograms of oriented gradients. Spatio-temporal pyramid matching [36] is used to generate a kernel for each visual word. Football plays are then classified into seven categories using multiple kernel learning. Hervieu et al. [19] use a hierarchical parallel semi-Markov model to classify different activity states in squash and handball, such as rallies, free throws and defence. The first layer describes the activity states, while the second layer consists of a parallel hidden Markov model for each feature representing the trajectories.

Perše et al. [49] represent team activity in basketball using team centroids to hierarchically classify situations with Gaussian Mixture Models. Thereafter, each situation is converted into a string, which is compared to templates for classification. Bialkowski et al. [5] use team centroids and occupancy maps to classify game situations in field hockey (corners, goals), emphasising the robustness of this representation to tracking noise. Grunz et al. [15] employ self-organising maps to identify long and short game initiations in soccer and Hirano et al. [20] use multi-scale comparison and rough clustering to cluster ball trajectories that lead to goals.

Direkoglu et al. [11] solve a special Poisson equation, in which the player positions determine the location of source terms. The derived distribution and its development over time defines a so-called region of interest used to describe the team behaviour. Wei et al. [61] use role models [43] and a Bilinear spatio-temporal basis model to represent team movement to cluster goal scoring opportunities in soccer. Bialkowski et al. [6] also use role models to automatically detect and compare the formations of soccer teams. Li et al. [39] learn a spatio-temporal driving force model to identify offence and defence players in football. Kim et al. [34] interpolate a dense motion field from player trajectories using thin-plate splines. This motion field is further investigated for points of convergence to predict where the game will evolve in short term.

From an application point of view, our approach is most comparable to Wei et al. [61] and Grunz et al. [15]. While Wei et al. focus on scoring opportunities and Grunz et al. study game initiations, we consider both situations in this study. Similar to Bialkowski et al. [5], our method proves robust to tracking noise.

## **3** Spatio-Temporal Convolution Kernels

#### 3.1 Representation

Multi-object trajectory analysis is concerned with a possibly varying number of moving objects  $\mathcal{O}_t$  in a set X, e.g.  $X = \mathbb{R}^2$ , over a finite period of time  $\mathcal{T} \subset \mathbb{N}$ . A multi-object trajectory is composed of snapshots of the object positions at different times. Depending on the context and application at hand, one of the following two formalisations of a snapshot is more appropriate.

**Definition 1 (Object-oriented Snapshot)** Assume the number of objects to be constant over time, i.e.  $\mathcal{O}_t = \mathcal{O} = \{o_1, ..., o_N\}$  for  $N \in \mathbb{N}$ . Then the object-oriented snapshot of all objects at time  $t \in \mathcal{T}$  is denoted by  $x_t \in X^N =: \mathcal{X}$ . We call  $\mathcal{X}$  the snapshot space. The position of a particular object  $o \in \mathcal{O}$  is denoted by  $x_t(o) \in X$ .

**Definition 2 (Group-oriented Snapshot)** Assume there is a constant number of groups  $K \in \mathbb{N}$ . Moreover, at every point in time each object can be

associated with exactly one of the groups  $\mathcal{G} = \{g_1, ..., g_K\}^{1}$ . Then the grouporiented snapshot of all objects at time  $t \in \mathcal{T}$  is denoted by

$$x_t \in \mathcal{P}(X)^K =: \mathcal{X}$$

We call  $\mathcal{X}$  the snapshot space. The positions of all objects of a particular group  $g \in \mathcal{G}$  are denoted by  $x_t(g) \in \mathcal{P}(X)$ . The group members of group g in snapshot  $x_t$  are denoted by  $O_{x_t}(g) \subset \mathcal{O}_t$ .

The implications of the two definitions are as follows. First, the object-oriented snapshot representation only allows a fixed number of objects, whereas the group-oriented representation is not limited in that respect. Second, in the group-oriented snapshot, objects inside a group are indistinguishable. On one hand, the property allows for permutations of objects but on the other hand it naturally also entails ambiguities.

Instead of an ordered sequence of positions or snapshots we use a set of time/position-pairs to represent trajectories. Thereby, time and order is explicitly represented as opposed to the more implicit sequence representation.

Definition 3 (Trajectory) A trajectory is defined as a finite subset

$$\tilde{P} = \{ (\tilde{t}_1, x_{\tilde{t}_1}), \dots, (\tilde{t}_n, x_{\tilde{t}_n}) \} \subset \mathcal{T} \times \mathcal{X},$$

such that  $\tilde{t}_i \neq \tilde{t}_j$  for  $i \neq j$ , i.e. the trajectory set  $\tilde{P}$  contains only one snapshot per point in time.

The set  $\pi_{\mathcal{T}}(\tilde{P}) = \{t \in \mathcal{T} : \exists (\tilde{s}, x_{\tilde{s}}) \in \tilde{P} \text{ s. t. } t = \tilde{s}\}$  contains all timestamps of the trajectory and is usually of the form  $\{K, K + 1, ..., K + L\}$  for some natural numbers K and L. When comparing trajectories it is insignificant at what absolute time the trajectories start. This gives rise to the following definition.

**Definition 4 (Time-Normalised Trajectory)** The time-normalised trajectory  $P \subset [0,1] \times \mathcal{X}$  corresponding to trajectory  $\tilde{P}$  is defined by normalising its time-scale to [0,1]. This corresponds to the trajectory P, given by

$$P = \{(t, x_t) : \exists (\tilde{t}, x_{\tilde{t}}) \in \tilde{P} \text{ s.t. } \tilde{t} = \mu + t(\max(\pi_{\mathcal{T}}(\tilde{P})) - \mu) \land x_{\tilde{t}} = x_t\},\$$

where  $\mu = \min(\pi_{\mathcal{T}}(\tilde{P})).$ 

In the remaining part of this study we refer to time-normalised trajectories simply by trajectories.

<sup>&</sup>lt;sup>1</sup> Note that the group membership of an object can change over time.

## 3.2 Problem Setting

One of the main advantages of kernel methods is the separation of algorithm and data. Following this paradigm, we focus on defining a kernel on the set of multi-object trajectories. Once a kernel has been defined, off-theshelf kernel machines can be applied to generate models, such as support vector machines [58], kernelised k-medoids [32], or spectral clustering [46]. The formal problem setting of this article is defined as follows. On the set of multiobject trajectories  $M \subset \mathcal{P}([0, 1] \times \mathcal{X})$  we aim to develop a similarity measure  $k : \mathcal{P}([0, 1] \times \mathcal{X}) \times \mathcal{P}([0, 1] \times \mathcal{X})$ , such that

- the absolute position as well as the shape of the trajectories is incorporated,
- (II) the measure is invariant to permutations of certain objects, i.e. for two trajectories  $P_1$ ,  $P_2$  with

$$P_2 = \{(t, x_t) : \exists \text{ permutation } \sigma \ \forall (s, y_s) \in P_1 \text{ s. t. } t = s \land x_t = \sigma(y_s)\}$$

it holds that  $k(P_1, P_2) = 1$ . In case of the group-oriented snapshot this already holds by definition if the permuted objects are members of the same group,

(III) the measure is invariant with respect to the speed of the movement. Since all trajectories have already been normalised to the same time scale, differences in speed are mainly reflected in the cardinality of the trajectory sets. So, for example, given two trajectories  $P_1$  and  $P_2$  with  $|P_1| = 2|P_2|$  and

$$P_2 = \{(t, x_t) : \exists (s, y_s) \in P_1 \text{ s. t. } t = 2s \land x_t = y_{2s}\}$$

it holds that  $k(P_1, P_2) \approx k(P_1, P_1)^2$ 

(IV) the similar movements of two sets of objects is recognised as such in the presence of deviations of single objects and outliers.

Moreover, the measure should have the following properties:

(V) Kernel Property, i.e.

$$k(P_1, P_2) = \langle \phi(P_1), \phi(P_2) \rangle_{\mathcal{F}}$$

for some, usually unknown, feature map  $\phi$  and inner product space  $\mathcal{F}$ 

- (VI) Broad applicability, i.e. few application specific parameters and no restrictions on the space X
- (VII) Computational efficiency

 $<sup>^2</sup>$  Note that the Definition of  $P_2$  is such that it corresponds to an object moving with twice the speed of the first object.

Note that properties (I) to (IV) formalise an intuitive notion of similarity. (I) says that shape and positions matters. (II) requires that if two similar object swap roles, it does not matter. (III) formalises that we do not care about differences in speed that much.<sup>3</sup>

Property (IV) demands robustness with respect to outlier trajectories. Further note that, for example, Dynamic Time Warping [4] meets condition (III) very well, but does not comply with (V), (VI) and (VII), since it is not a kernel and only applicable if the underlying set is a metric space. Moreover, it is computationally expensive. On the other hand, the Hausdorff distance [17] satisfies (I), (III) and (VII), but it does not satisfy (II), (V) and (VI), since it is only applicable to metric spaces and sensitive to permutations. In addition, it is not a kernel. A Gaussian RBF kernel on the full vector of positions meets conditions (I) (restricted), (V) and (VII), but is not applicable to sequences of different lengths. Also, it does not comply with (II),(III) and (VI), since it is highly sensitive to variations in speed and permutations and is restricted to metric spaces.

#### 3.3 Spatio-temporal Convolution Kernels for Multi-trajectories

In this section we develop a kernel on the space of (time-normalised) multitrajectories  $\mathcal{P}([0,1] \times \mathcal{X})$ . Each of those trajectories consists of a set of snapshots associated with a relative time. The general idea is to perform a pairwise comparison of the snapshots in the two sets. Therefore, we first need a way to compare snapshots and, second, we need to know which snapshots of the two trajectory sets to compare with each other. For the latter dynamic time warping (DTW) [4] seems to be a good choice, since it aligns the snapshots optimally in terms of similarity. Unfortunately, the obtained kernel is not positive definite, i.e. it does not correspond to an inner product in some Hilbert Space. Although there is anecdotal evidence that learning with indefinite kernels can lead to good results in some applications (e.g. [47]), theory only supports the use of positive definite kernels. For many kernel machines there are error bounds and convergence criteria that can be straightforwardly applied to positive definite kernels but that do not hold for indefinite kernels [7,41,56].

Therefore, we propose a weighted comparison between every snapshot of the first trajectory and every snapshot of the second one where the weights depend on the offset in relative time. Formally, this is done using an Rconvolution kernel [18] on the two sets representing the trajectories. Convolution kernels are a general class of kernels on structured objects  $x, y \in$ X.<sup>4</sup> The idea is to compare instances x and y by comparing their parts

<sup>&</sup>lt;sup>3</sup> Note that the property can be easily adapted to explicitly include speed by adding an extra coordinate. The position of each object is replaced by a position-speed-pair  $\left(x, \frac{\mathrm{d}x}{\mathrm{d}t}\right)_t$ , see [27] for details.

<sup>&</sup>lt;sup>4</sup> Originally, convolution kernels are defined on arbitrary sets X.

 $(x_1, ..., x_D), (y_1, ..., y_D) \in X_1 \times ... \times X_D$ . Thus, a relation function R is needed to express that something is a *part* of some structure.

**Definition 5 (Relation)** Let  $X, X_1, ..., X_D$  be arbitrary sets. Then a relation

$$R: X \times X_1 \times \dots \times X_D \to \{0, 1\}$$

is an arbitrary boolean function that returns 1 if and only if  $(x_1, ..., x_d)$  are parts of x. The set of parts of  $x \in X$  under relation R is denoted by

$$R^{-1}(x) = \{(x_1, ..., x_d) \in X_1 \times ... \times X_D : R(x, x_1, ..., x_D) = 1\}$$

and R is called **finite** if  $R^{-1}(x)$  is finite for every  $x \in X$ .

**Definition 6** (*R*-Convolution Kernel) Let  $X, X_1, ..., X_D$  be arbitrary sets. Let  $x, y \in X$  and  $R: X \times X_1 \times ... \times X_D \to \{0, 1\}$  be a finite relation. Moreover, let  $k_1, ..., k_D$  be kernels on  $X_1, ..., X_D$ . Then the *R*-convolution kernel on X is defined by

$$k(x,y) = \sum_{\substack{(x_1,\dots,x_D)\in R^{-1}(x), \ d=1\\(y_1,\dots,y_D)\in R^{-1}(y)}} \prod_{d=1}^D k_d(x_d, y_d)$$

The following theorem shows that an *R*-convolution kernel is indeed a (positive-definite) kernel.

**Theorem 1** Let  $X, X_1, ..., X_D$  be arbitrary sets. Let R be a finite relation and let  $k_1, ..., k_D$  be kernels on  $X_1, ..., X_D$ . Then the R-convolution kernel k given by Definition 6 is a kernel.

*Proof* For the proof we refer to [18] Theorem 1 and Lemma 1, which are essentially more involved applications of closure properties of kernels.  $\Box$ 

In our case the structure is a multi-object trajectory and its *parts* are the snapshots, the times of the snapshots and the length of the trajectory. The relation between the structure and its elementary components is given by

$$R: \underbrace{\mathbb{N}}_{\text{Length of Traj.}} \times \underbrace{[0,1]}_{\text{Time of Snapshot}} \times \underbrace{\mathcal{X}}_{\text{Snapshot}} \times \underbrace{\mathcal{P}([0,1] \times \mathcal{X})}_{\text{Trajectory}} \to \{0,1\}$$

To extract the snapshots with their associated times and the length of the trajectory from the trajectory, R needs to be defined as follows:

$$R(n,t,x,P) = \begin{cases} 1 & \text{if } |P| = n \land \exists (s,y_s) \in P : (t,x) = (s,y_s) \\ 0 & \text{otherwise.} \end{cases}$$
(1)

The first condition in Equation 1 ensures that the given trajectory has the given length, while the second condition guarantees the occurrence of the given

snapshot at the given time in the given trajectory. With this relation and the kernels to compare the *parts* the resulting convolution kernel is given by

$$k(P,Q) = \sum_{\substack{(n,t,x)\in R^{-1}(P), \\ (m,s,y)\in R^{-1}(Q)}} k_{\mathbb{N}}(n,m) \cdot k_{[0,1]}(t,s) \cdot k_{\mathcal{X}}(x,y)$$
(2)

with  $R^{-1}(P) = \{(n, t, x) : R(n, t, x, P) = 1\}$ . The term  $k_{\mathbb{N}}$  accounts for differences in length of the trajectories by normalisation, i.e.  $k_{\mathbb{N}} = \frac{1}{nm}$ . Note that  $k_{\mathbb{N}}$  is a kernel on  $\mathbb{N}$ , since it corresponds to the standard inner product under the feature map  $\phi : \mathbb{N} \to \mathbb{N}$  given by  $\phi : n \mapsto \frac{1}{n}$ . Finally, the *R*-convolution kernel simplifies to

$$k(P,Q) = \frac{1}{|P||Q|} \sum_{(t,x_t) \in P, (s,y_s) \in Q} k_{[0,1]}(t,s) \cdot k_{\mathcal{X}}(x_t,y_s).$$
(3)

**Theorem 2** The spatio-temporal convolution kernel (Equation 3) is a kernel if the temporal kernel  $k_{[0,1]}$  and the spatial kernel  $k_{\chi}$  are kernels in that sense.

Proof By Theorem 1 we need to show that R is finite and the component kernels  $k_{\mathbb{N}}$  and  $k_{[0,1]}$ ,  $k_{\mathcal{X}}$  are indeed kernels. First, for all  $P \in \mathcal{P}([0,1] \times \mathcal{X})$  it holds that  $|R^{-1}(P)| = |P| < \infty$  by Definition 3, so R is finite. Second,  $k_{[0,1]}$ and  $k_{\mathcal{X}}$  are kernels by assumption and we have just shown that  $k_{\mathbb{N}}$  is also a kernel. Hence, the spatio-temporal convolution kernel is a kernel.

As indicated, Equation 3 can be interpreted such that all snapshots of the two trajectories are compared with each other, but weighted by their offset in time. Thereby snapshots which occur at different relative times have a low contribution to the overall similarity, while snapshots at the same relative time have a high contribution. This is why  $k_{[0,1]}$  is sometimes also referred to as weight function. The definition of k in Equation 3 leaves two degrees of freedom:

- Spatial kernel  $k_{\mathcal{X}}$ : The choice of the snapshot kernel determines which snapshots are similar.
- Temporal kernel  $k_{[0,1]}$ : The choice of the temporal kernel determines the way in which the snapshots of two sequences are combined, and thus the importance of ordering and speed.

In the following subsections we develop and compare different spatial and temporal kernels.

## 3.4 Spatial Kernels

A spatial kernel compares two snapshots in X. Corresponding to the two definitions of the snapshot in Definition 1 (object-oriented) and Definition 2 (group-oriented), two types of kernels are introduced here as well.

Name	X	Kernel	Para.
Uniform	$\mathbb{R}^{n}$	$k_X(x,y) = \mathbb{I}_{\{z \in \mathbb{R}^n : \ z-x\ _2 < w\}}(y)$	w
Triangular	$\mathbb{R}^{n}$	$k_X(x,y) = \left(1 - \frac{\ x-y\ _2}{w}\right) \mathbb{I}_{\{z \in \mathbb{R}^n : \ z-x\ _2 < w\}}(y)$	w
Polynomial	$\mathbb{R}^{n}$	$k_X(x,y) = (\langle x,y \rangle + R)^d$	R,d
Gaussian	$\mathbb{R}^{n}$	$k_X(x,y) = \exp\left(-\frac{1}{2\sigma^2}  x-y  _2^2\right)$	$\sigma^2$
Matching Kernel	finite $X$	$k_X(x,y) = \mathbb{I}_{\{x\}}(y)$	-

Table 1 Elementary Spatial Kernels

*Object-Wise Comparison* These kernels correspond to the object-oriented snapshot by simply comparing the positions of the objects  $\mathcal{O}$  and summing up their similarity:

$$k_{\mathcal{X}}(x_t, y_t) = \frac{1}{|\mathcal{O}|} \sum_{o \in \mathcal{O}} k_X(x_t(o), y_t(o)).$$

$$\tag{4}$$

Note that Equation 4 is a kernel, since kernels are closed under direct sum and multiplication by a positive constant (see [54] Proposition 3.22).<sup>5</sup> Since kernels are also closed under direct product, technically a product could have been used instead of the sum in Equation 4. However, a product of kernels leads to vanishing similarities if only one object is dissimilar to its counterpart. This is counterintuitive as two snapshots are the more similar the more objects match well. It is thus an inherently additive relationship. For the same reason a product of kernels is more vulnerable to noise and outliers and generally less robust than a sum of kernels. The elementary kernel  $k_X$  can be any kernel comparing two positions in X. Table 1 lists common kernels for finite as well as continuous snapshot spaces.

Kernels as in Equation 4 have two major shortcomings. First, they penalise permutations of objects. For example, two snapshots of two objects with swapped positions will have low similarity, although in terms of the group motion they are alike. One way to address permutations is to explicitly maximise the similarity of the two snapshots with respect to all possible permutations of objects. Due to the high computational costs of considering all possible permutations<sup>6</sup>, this is infeasible. In addition, the bandwidth parameter  $\sigma$  that controls the interval of high sensitivity, i.e. large values of |dk/d(||x - y||)|, of the kernel has to be known beforehand. This is critical since one usually does not know on which scale significant deviations appear. Both issues are addressed by the group-wise comparison of objects.

Group-Wise Comparison In case of the group-oriented snapshot there is a partition of  $\mathcal{O}_t$  into K sets  $O_{x_t}(g_1), O_{x_t}(g_2), .., O_{x_t}(g_K)$ . Instead of comparing

<sup>&</sup>lt;sup>5</sup> However, if the number of objects was not constant,  $k_{\mathcal{X}}$  would not be a kernel.

<sup>&</sup>lt;sup>6</sup> For example, 10 objects lead to about  $3.6 \cdot 10^6$  permutations.

N objects, K groups of objects are compared, which leads to the following definition of a spatial kernel

$$k_{\mathcal{X}}(x_t, y_t) = \frac{1}{K} \sum_{g \in \mathcal{G}} k_G(x_t(g), y_t(g))$$
(5)

with kernel  $k_G$  applied on sets of positions. For the same reasons as in the previous paragraph, a sum of kernels is preferable to a product of kernels. Definition 2 allows for zero group members in a snapshot, i.e.  $x_t(g) = \emptyset$ . For all kernels defined below this special case is dealt with as follows:

$$k_G(x,y) = \begin{cases} 1 & \text{if } x = y = \emptyset \\ 0 & otherwise. \end{cases}$$

The definition of  $k_G$  is very much dependent on the underlying set of positions X. There are two main classes of sets X. First, if X is a metric space, i.e. there exists a distance measure between all positions. Second, X is a finite set, in which no meaningful distance between positions can be defined.<sup>7</sup> In the case of a metric space we will subsequently only consider  $X = \mathbb{R}^2$  with standard Euclidean distance as it is the most common in applications and can be directly extended to  $X = \mathbb{R}^N$ .

Euclidean Space:  $X = \mathbb{R}^2$  In the case of  $X = \mathbb{R}^2$  a straightforward approach is to use one of the kernels from Table 1 to compare the group centroids

$$\mu_{x_t}(g) = \frac{1}{|O_{x_t}(g)|} \sum_{o \in O_{x_t}(g)} x_t(o) \quad \text{and} \quad \mu_{y_t}(g) = \frac{1}{|O_{y_t}(g)|} \sum_{o \in O_{y_t}(g)} y_t(o)$$

which leads to

$$k_G(x_t(g), y_t(g)) = k_X(\mu_{x_t}(g), \mu_{y_t}(g)).$$

Besides the need of defining the width parameter  $\sigma$ , disregarding the distribution of the objects around their centroids constitutes a disadvantage. To remedy both deficiencies, two Gaussian distributions are fitted to  $x_t(g)$ , and  $y_t(g)$  respectively that we will denote as  $f_{x_t}(g)$  and  $f_{y_t}(g)$ . These distributions are defined by their means  $\mu_{x_t}(g)$  and  $\mu_{y_t}(g)$  as defined above and their covariance matrices defined by

$$\Sigma_{x_t}(g) = \frac{1}{|O_{x_t}(g)| - 1} \sum_{o \in O_{x_t}(g)} (x_t(o) - \mu_{x_t}(g)) (x_t(o) - \mu_{x_t}(g))^T$$
  
and  $\Sigma_{y_t}(g) = \frac{1}{|O_{y_t}(g)| - 1} \sum_{o \in O_{y_t}(g)} (y_t(o) - \mu_{y_t}(g)) (y_t(o) - \mu_{y_t}(g))^T.$ 

 $<sup>^{7}</sup>$  The case of an infinite set X without a metric can be reduced to the finite case by only considering the positions that are attained by some objects at some time. These are only finitely many.

If the covariance matrix is ill-conditioned or singular<sup>8</sup>, a simple shrinking scheme with shrinkage parameter  $\alpha$  can be applied to achieve non-singularity:

$$\Sigma = (1 - \alpha) \cdot \Sigma + \alpha \cdot \frac{\operatorname{Tr}(\Sigma)}{2} \mathbb{I}_2,$$

where  $\mathbb{I}_2$  is the two-dimensional identity matrix. There exist different strategies [37,10] to choose an optimal value for  $\alpha$ , but for our purposes it suffices to deploy a constant 0.1. In the case of  $\text{Tr}(\Sigma) = 0$  the following scheme is used

$$\Sigma = (1 - \alpha) \cdot \Sigma + \alpha \cdot \sigma_{\mathrm{MIN}}^2 \cdot \mathbb{I}_2 = \alpha \cdot \sigma_{\mathrm{MIN}}^2 \cdot \mathbb{I}_2,$$

where  $\sigma_{\text{MIN}}^2$  is an application specific parameter. Strategies for choosing  $\sigma_{\text{MIN}}^2$  are discussed in Section 4. Reasons for the matrix to be singular are  $O_t(g) \leq 2$  or strong collinearity of the samples. The two Gaussian distributions  $f_{x_t}(g)$  and  $f_{y_t}(g)$  are then compared using a probability product kernel [25], i.e.

$$k_G(x_t(g), y_t(g)) = k^{\text{prod}}(f_{x_t}(g), f_{y_t}(g)).$$

A probability product kernel compares two probability distributions, i.e. their density functions p and p', defined on a common probability space  $(\Omega, \mathcal{A}, \mu)$ . It is defined as follows.

**Definition 7 (Probability Product Kernel)** Let p and p' be probability density functions defined on the same probability space  $(\Omega, \mathcal{A}, \mu)$ . Assume that  $p^{\rho}, p'^{\rho} \in L^2(\Omega)$  for  $\rho \in \mathbb{R}^+$ . Then the probability product kernel (PPK)  $k^{\text{prod}} : L^2(\Omega) \times L^2(\Omega) \to \mathbb{R}$  is given by

$$k^{\mathrm{prod}}(p,p') = \int_{\Omega} p(\omega)^{\rho} p'(\omega)^{\rho} \,\mathrm{d}\mu(\omega).$$

**Lemma 1**  $k^{prod}$  is a (positive-definite) kernel.

*Proof* Consider the feature map  $\phi: L^2(\Omega) \to L^2(\Omega)$  given by

$$\phi(p) = p^{\rho},$$

which is well-defined because of above  $L^2$  property. Then

$$k^{\operatorname{prod}}(p,p') = \langle \phi(p), \phi(p') \rangle_{L^2(\Omega)}$$

and  $k^{\text{prod}}$  is a kernel corresponding to the scalar product in  $L^2(\Omega)$ 

$$\int_{\Omega} f(\omega) \cdot g(\omega) \,\mathrm{d}\mu(\omega).$$

E.	-	7	
L		L	
L			

 $<sup>^{8}\,</sup>$  In this study a threshold on the covariance matrix's condition number of 30 is used.

In this study we will use  $\rho = 1/2$ , which has the important property that  $\text{Im}(k_{Prod}) = [0,1]$  and  $k^{\text{prod}}(p,p) = 1$ . For two two-dimensional Gaussian distributions and  $\rho = 1/2$  the product probability kernel permits the following closed form:

$$k_{G}(x_{t}, y_{t}) = \int_{\mathbb{R}^{2}} \left( f_{x_{t}}(z) f_{y_{t}}(z) \right)^{1/2} dz = 2 \cdot |\Sigma^{*}|^{\frac{1}{2}} |\Sigma_{x_{t}}|^{-\frac{1}{4}} |\Sigma_{y_{t}}|^{-\frac{1}{4}} \cdot \exp\left( -\frac{1}{4} \left( \mu_{x_{t}}^{T} \Sigma_{x_{t}}^{-1} \mu_{x_{t}} + \mu_{y_{t}}^{T} \Sigma_{y_{t}}^{-1} \mu_{y_{t}} - \mu^{*T} \Sigma^{*} \mu^{*} \right) \right)$$
(6)

with  $\Sigma^* = (\Sigma_{x_t}^{-1} + \Sigma_{y_t}^{-1})^{-1}$  and  $\mu^* = \Sigma_{x_t}^{-1} \mu_{x_t} + \Sigma_{y_t}^{-1} \mu_{y_t}$ .<sup>9</sup> Besides its sensibility to the distribution of the objects inside the group  $k^{\text{prod}}$  has the advantage of not having to choose a width parameter. The kernel essentially adapts to the scale of the group.

Finite Case:  $|X| = L < \infty$  For finite  $X = \{z_1, ..., z_L\}$ , we can generalise to the group-oriented snapshot by comparison of the number of objects of each group at each possible location in X:

$$k_G(x_t(g), y_t(g)) = \frac{1}{L} \sum_{z \in X} \mathbb{I}_{\{n_z(x_t(g))\}}(n_z(y_t(g)))$$
(7)

with the indicator function  $\mathbb{I}_A(x) = 1$  if  $x \in A$  and 0 otherwise and  $n_z(x_t(g)) = |\{x_t(o) = z \ o \in O_{x_t}(g)\}|$  denotes the number of objects of group  $g \in \mathcal{G}$  in position  $z \in X$  at time  $t \in [0, 1]$ .

**Lemma 2**  $k_G$  as defined in Equation 7 is a kernel.

*Proof* Let  $l^2(0,1)$  be the Hilbert space of 0/1-sequences equipped with the standard scalar product. Consider  $\phi: X \to l^2(0,1)^L$  with

$$\phi(x_t(g))_{m,l} = \begin{cases} 1 & \text{if m objects in position l} \\ 0 & \text{otherwise} \end{cases}$$

for l = 1, ..., L and m = 0, 1, 2, ... Using this feature mapping,  $k_G$  can be rewritten to

$$k_G(x_t(g), y_t(g)) = \langle \phi(x_t(g)), \phi(y_t(g)) \rangle$$

showing that it is indeed a kernel. Note that that the inner product of a product space is defined by the sum of the inner products of its components.  $\Box$ 

If the number of possible locations is much higher than the number of objects, the above kernel (Equation 7) is inappropriate, as it will return similarities close to one for every pair of snapshots because in both snapshots there are no objects at almost all locations. Alternatively, the idea of a product probability kernel can also be applied to the finite setting by fitting a multinomial

<sup>&</sup>lt;sup>9</sup> In order to simplify the notation, group index g has been omitted.

distribution to the positions of the group and comparing the two resulting distributions  $p_{x_t}(g)$  and  $p_{y_t}(g)$  that are identified by their outcome probabilities  $(p_{x_t,1}(g), ..., p_{x_t,L}(g))$ . For  $\rho = 1/2$  the probability product kernel on the multinomial distributions is derived as follows<sup>10</sup>,

$$\begin{aligned} k^{\text{prod}}(x_t, y_t) &= \int_{\mathbb{N}^L} \sqrt{p_{x_t}(n_1, \dots, n_L) p_{y_t}(n_1, \dots, n_L)} \, \mathbf{d}(n_1, \dots, n_L) \\ &= \sum_{n_1 \in N} \dots \sum_{n_L \in N} \sqrt{p_{x_t}(n_1, \dots, n_L) p_{y_t}(n_1, \dots, n_L)} \\ &= \sum_{n_1 \in N} \dots \sum_{n_L \in N} \binom{n_1 + \dots + n_L}{n_1, \dots, n_L} \left( (p_{x_t, 1} \cdot p_{y_t, 1})^{n_1} \cdot \dots \cdot (p_{x_t, L} \cdot p_{y_t, L})^{n_L} \right)^{\frac{1}{2}} \end{aligned}$$

with the maximum-likelihood estimation for the outcome probabilities of the multinomial distributions given by

$$p_{x_t,l}(g) = \frac{n_{z_l}(x_t(g))}{|O_{x_t}(g)|}.$$

## 3.5 Temporal Kernels

The temporal kernel  $k_{[0,1]}$  is the simpler component of spatio-temporal convolution kernel because the underlying set is fixed to the one-dimensional interval [0, 1]. We briefly discuss possible options for the temporal kernel and their implications.

The simplest temporal kernel is a constant kernel  $k_{[0,1]}(t,s) = 1$  with the consequence that the spatio-temporal convolution kernel (Equation 3) collapses to a set kernel on the two sets of snapshots, thus ignoring order. A uniform or interval kernel given by

$$k_{[0,1]}(t,s) = \mathbb{I}_{\{u \in \mathbb{R} : |t-u| < w\}}(s)$$

could be used to compare every snapshot of the first trajectory to just those snapshots of the second trajectory which are close in (relative) time. The choice of w determines how close snapshots have to be to be taken into consideration. Finally, Gaussian kernels of the form

$$k_{[0,1]}(t,s) = \exp\left(-\frac{1}{2\sigma^2} \|t-s\|_2^2\right)$$

compare every snapshot of the first trajectory to every snapshot of the second trajectory and gives more importance to closer events (in relative time). Depending on the application at-hand, there may be many more options to define temporal kernels.

 $<sup>^{10}\,</sup>$  The group index g has again been omitted for reasons of readability.

	gorithing i filo i offer i offernaal ripproximation
1:	<b>function</b> ASTCKII(P,Q,L) $\triangleright$ Input: Two Trajectories P and Q and the max. length of any trajectory L. Output: Vector of step by step approximation of $k(P,Q)$
	of any trajectory $L$ , Output. Vector of step-by-step approximation of $K(1,Q)$
2:	tempMatrix = getTemporalKernelMatrix(P,Q);
3:	orderedEntries = sort(tempMatrix)
4:	$\mathrm{k}(0)=0;$
5:	for i=1 to L do
6:	$\mathrm{k}(\mathrm{i}) = \mathrm{k}(\mathrm{i}\text{-}1)\cdot\mathrm{round}((\mathrm{i}\text{-}1)/\mathrm{L}\cdot \mathrm{P} \cdot \mathrm{Q} );$
7:	for j = $ROUND((i-1)/L  P   Q )+1$ to $ROUND(i/L  P   Q )$ do
8:	k(i) += spatialKernel(P(orderedEntries(j,1)) Q(orderedEntries(j,2)));
9:	end for
10:	$\mathrm{k(i)} = \mathrm{k(i)}/\mathrm{round}(\mathrm{i}/\mathrm{L}\cdot \mathrm{P} \cdot \mathrm{Q} )$
11:	end for
12:	return k;
13:	end function

# Algorithm 1 ASTCK: Percental Approximation

3.6 Approximation Techniques

To compute the Gram matrix of a dataset of N trajectories,  $O(N^2L^2)$  spatial as well as temporal kernels need to be evaluated with L being the maximal length of a sequence. Naturally, the evaluation of the spatial kernels will dominate the evaluation of the temporal kernel mainly for reasons of dimensionality as well as complexity of the kernel itself.

Generally, there exist two ways to speed up the algorithm. First, the number of similarities which have to be computed is reduced to shrink the factor  $O(N^2)$ . Afterwards the missing entries in the Gram matrix are reconstructed. Second, the computation of each similarity is accelerated to reduce the term  $O(L^2)$ . Regarding the first approach, there exist methods to reconstruct incomplete Gram matrices, most notably the Nyström method (see [12,62]). These methods have the disadvantage of usually needing a constant fraction of all entries to perform well [60]. It follows that they do not improve the asymptotic complexity of the computation. We focus on the second approach, since it is specific to spatio-temporal convolution kernels and potentially improves the asymptotic complexity.

As stated before, the computation of the temporal kernel is significantly faster than the one of the spatial kernel. Thus, a natural way to approximate these kernels is to first evaluate all temporal kernels and then to only evaluate those spatial kernels which have the highest contribution. Depending on the lengths of the two trajectories at-hand the number of overall kernel evaluations varies between pairs of trajectories. Therefore, when approximating the kernel one has to take care that all entries of the Gram matrix are approximated in similar speed to avoid distortions. To this end, we propose a percental approximation algorithm (Algorithm 1, Approximated STCK (ASTCK)), which in addition to the two trajectories takes the maximal length of all trajectories L as an input. The algorithm first evaluates all temporal kernels. Subsequently, in each step  $\frac{1}{L}$  spatial kernels are computed starting with the most relevant based on the evaluation of the temporal kernels. The order of evaluation is

Table 2 Evaluation Order of ASTCK with L = 12 and |P| = |Q| = 6

		Timesteps $P$						
	$\mathbf{Step}$	0	.2	.4	.6	.8	1	
Timesteps $Q$	0	1	3	6	9	11	12	
	.2	4	1	3	6	9	11	
	.4	7	5	1	3	7	9	
	.6	10	8	5	<b>2</b>	4	7	
	.8	11	10	8	5	<b>2</b>	4	
	1	12	12	10	8	6	<b>2</b>	

illustrated in Table 2. This algorithm leads to a complexity of

$$\mathcal{O}(N^2 \cdot L^2 \cdot C_T + N^2 \cdot L \cdot I \cdot C_S),$$

where  $C_T$  is the complexity of evaluating the temporal kernel,  $C_S$  is the complexity of evaluating the spatial kernel and I is the number of iterations of the approximation algorithm. The parameter  $1 \leq I \leq L$  needs to be chosen by the user. In practice, small values of I often lead to highly accurate results.

#### 3.7 Online Application

STCKs are particularly well-suited for online analyses such as real-time computations where the objects of interest are still in motion. When a new measurement is added to a trajectory, Equation 3 can be efficiently updated, since all previous spatial kernel evaluations remain constant in the sum. That is, only the value of the temporal kernel needs to be computed which is however usually inexpensive.

#### **4 Empirical Evaluation**

In this section we empirically compare our spatio-temporal convolution kernels to baseline approaches using artificial and real data sets. We focus on clustering tasks and k-medoids [32] as the underlying learning algorithm. The temporal kernel is always a Gaussian kernel that we combine with three different spatial kernels: an object-wise Gaussian RBF kernel as spatial kernel (STCK<sub>1-1</sub>), a Gaussian RBF kernel on the group means (STCK<sub>Mean</sub>), and a probability product kernel on the fitted Gaussian distributions (STCK<sub>Dist</sub>).

For the latter, the group and application specific parameter  $\sigma_{\text{MIN}}^2$  which is sometimes needed to restore non-singularities of the covariance matrix (see Section 3.4), is set to the average distance between two objects of the group. If the group only consists of one object, it is equal to the average distance between two objects of the group, that is most similar. The width parameter  $\sigma_T$  of the temporal Gaussian kernel is set to 0.5 to balance invariance to speed differences and sensitivity to direction. The width parameter  $\sigma_S$  of the spatial Gaussian RBF kernel is set to 0.2.

We deploy three baselines. First, Junejo et al. [28] is straightforwardly extended to the multi-object scenario, i.e. the use of Hausdorff distance on the set of positions of the trajectories. Instead of the hierarchical clustering employed in [13], kernelised k-medoids is used. The second baseline is inspired by Grimson et al. [59]. We use a bag-of-positions as well as a bag-of-directions representation for the trajectories of each group. To keep the setup simple, we use a multinomial mixture model (MNMM) and expectation maximisation for clustering instead of a semantic topic model like dual-HDP [59]. Third, we also compare our method to dynamic time warping with a product probability kernel (DTW<sub>dist</sub>) serving as local distance measure. This method applies the product probability kernel to the fitted Gaussian distributions. The number of clusters is determined using the silhouette measure [52], Hartigan index [16] as well as next-neighbour consistency for all methods.

In the next section, we measure the alignment between predicted and ground-truth clusterings using artificially generated data. The alignment between two groupings is captured by the Rand Index [51]; however, as two random clusterings may return a non-zero Rand Index by chance, we resort to the Adjusted Rand Index [22], given by

$$AR(S,T) = \frac{R(S,T) - \mathbb{E}[R(S,T)]}{1 - \mathbb{E}[R(S,T)]}$$

where R denotes the Rand Index and S and T are clusterings to compare.

#### 4.1 Artificial Data

Recall that our spatio-temporal convolution kernels possess basic properties<sup>11</sup> that are not shared with the baseline methods. Most notably these are invariance to permutations of the objects and to differences in speed as well as sensitivity to the spatial distribution of the objects and the direction of the movement.

In this section these properties are experimentally confirmed using two toy data sets consisting of artificially generated multi-object trajectories. The data is generated to cover a broad range of trajectories present in real-world applications. The first setting covers linear movements and each trajectory consists of five objects performing a linear movement as shown in Figure 1 (top). To find the correct clusters, the respective kernels need to distinguish trajectories based on the direction of the movement or based on the distribution of the objects, while the objects' centroid is the same in both trajectories.

The second setting deals with circular movements and each trajectory consists of four objects moving in circles of different radii, see Figure 1 (bottom). The construction of the second set enables us to evaluate the ability of the

<sup>&</sup>lt;sup>11</sup> The following properties refer to the  $STCK_{Dist}$ 



Fig. 1 Exemplary artificial multi-trajectories. For both settings ten multi-trajectories for each cluster are shown. Colours indicate cluster membership. Moreover, two levels of noise are depicted. For instance, *linear 0.025* show linear movements with added noise in the interval [-0.025, 0.025].

kernels to identify the correct clusters when these only differ in the direction of the movement but a spatial separation between the clusters is not possible. Compared to linear movements, the circular task is more difficult, as the directions of only some objects differ between the clusters. The data generation is described in detail in Appendix A.

To evaluate the sensitivity of the methods to permutations, the objects' ordering inside each multi-object trajectory is permuted randomly in both data sets. To assess the sensitivity to changes in speed, we flip a coin for each trajectory in both test sets. With 50% probability only every second position is retained in the trajectory. The resulting trajectory corresponds to a movement with twice the speed of the original trajectory. Finally, we add uniformly distributed noise in the range  $[-\epsilon, \epsilon]$  for  $\epsilon \in \{0.005, 0.025, 0.05, 0.1, 0.25\}$  with zero mean to every position in each multi-object trajectory in both test sets.

For both data sets, we generate 200 multi-trajectories (50 per cluster) for all five levels of noise. It is important to assess the methods' sensitivity to noise for two reasons. First, trajectories are generally highly variable and exact matches



Fig. 2 Cluster performance of the STCK compared to three baselines



Fig. 3 Comparison of ASTCK II to the baseline methods

are very rare. Second, due to frequent tracking errors, real-life trajectories are also subject to a significant level of noise.

In this experiment, we focus on  $\text{STCK}_{\text{Dist}}$  as it leads to more accurate clusters than  $\text{STCK}_{1-1}$  and  $\text{STCK}_{\text{Mean}}$  throughout all ranges of noise. Figure 2 compares  $\text{STCK}_{\text{Dist}}$  to the three baselines. Our method outperforms dynamic time warping on both test sets and gives the most accurate clusterings. The multinomial mixture model performs reasonably well on linear movements but leads to inappropriate clusterings for circular ones. The poor results on the second data set are due to the absence of ordering information in the bag-of-directions representation when the objects move in a full circle. The Hausdorff distance-based Junejo et al. leads to generally inaccurate clusterings. We credit



Fig. 4 Runtime: number of trajectories (top), length of trajectories (bottom, left), and number of objects (bottom, right).

this finding to its sensitivity to permutations as well as negligence of ordering information.

We now evaluate the proposed approximation technique on both artificial data sets. The algorithm performs up to eleven iterations as the longest trajectory consists of eleven snapshots. Note that, at the end of the last iteration, the exact Gram matrix is obtained. The percental approximation proves accurate; the approximation quickly leads to the correct clustering after two iterations. Figure 3 shows that ASTCK also achieves higher consistencies with the correct clustering compared to all baseline methods from the second iteration onwards. Moreover, starting from the second iteration onwards it is always at least as good as the exact method. The depicted results correspond to a noise level of 0.005 but equivalent results hold for all other noise ratios.

The runtime of multi-object trajectory clustering is governed by the number of trajectories, the length of the trajectories and the number of objects. Depending on the application, usually one or two dimensions are dominating. It is thus important to evaluate the runtime for each of the three dimensions. The theoretical runtime is  $O(N^2 \cdot L^2)$  spatial as well as temporal kernel evaluations, where N is the number of trajectories and L is the maximum length of the trajectories. To confirm these findings experimentally, we generate random trajectory sets in  $[0, 1]^2$  and vary the number of trajectories per set, the length of the trajectories and the number of objects per trajectory. The results are depicted in Figure 4.

The figure on top shows the results for varying numbers of trajectories. Except for multinomial mixture models that perform linearly in the number of instances, all methods exhibit a quadratic complexity, which is simply due to the fact that all pairwise similarities are computed, i.e. N(N + 1)/2. Varying the length of the trajectories in Figure 4 (bottom, left) shows that all STCKs as well as DTW exhibit quadratic complexities which is in line with the theoretical runtime. However, the percental approximation algorithm (two iterations) significantly improves the runtime of STCKs and is comparable to that of Junejo et al. and multinomial mixture models.

Finally, Figure 4 (bottom, right) varies the number of objects. The results are virtually constant time complexities for all methods except  $STCK_{1-1}$  and MNMM, which exhibit linear complexities. The observation is in line with our expectations as the Gaussian RBF kernel compares each object with its counterpart. The deviations in the case of  $STCK_{Dist}$  for a small number of objects (less or equal two objects) is due to the additional time needed by the shrinking schemes to restore non-singularity of the covariance matrices.

## 4.2 Real-world Data

We now evaluate spatio-temporal convolution kernels on real world data using positional data streams of ten soccer games of the German Bundesliga. The goal is to identify movement patterns by analysing the tracking data.

The tracking data is captured by the VIS.TRACK [23] tracking system during five games of Bundesliga Team A and five games of Bundesliga Team B from the 2011/12 Bundesliga season.<sup>12</sup> VIS.TRACK is a video based tracking system consisting of several cameras in the soccer stadium. It determines the positions of the players, ball and referees at 25 frames per second, which amounts to roughly 135,000 positions per object and match and a total of 31,000,000 positions. Additionally, a marker indicates the status of the ball (in-play, stoppage) and the possession of the ball. The range of the x- (parallel to sidelines) and y-coordinate (parallel to endlines) is [-1, 1], whereas values with an absolute value greater than one occur if the ball is out of bounds. The data stream is preprocessed so that positions of the second half are mirrored to account for the changeover at half time and the frame numbers of the second half are changed to succeed those of the first half. Additionally, the playing direction is determined and normalised, so that the team of interest always plays from left to right. Subsequently, we extract two types of sequences: game initiations and scoring opportunities.

Game initiations (GI) begin with the goal keeper passing the ball and end with the team loosing possession, a stoppage, the ball being in the attacking third of the field or the start of the next game initiation as defined above.

 $<sup>^{12}\,</sup>$  The team names must not be disclosed.



Fig. 5 Pairwise adjusted Rand indeces.



Fig. 6 Cluster consistencies and compactness.

Sequences shorter than length 12 are excluded. Scoring opportunities (SO) terminate when the ball is carried into a predefined zone of danger, usually defined as the attacking quarter. Scoring opportunities begin at the time of the last stoppage or win of the ball and last until the ball reaches the attacking quarter of the field  $[0.5, 1] \times [-1, 1]$ . Again, sequences with a length below 12 are discarded.

For every sequence there are 23 possible trajectories (ball, 22 players, no referees) to include into the analysis. Since the opposing team changes from game to game, we will restrict the analysis to twelve objects, namely the ball and the players of Team A, and Team B respectively. In the following experiments we consider the ball as one group in the sense of Definition 2. We further include the back four (game initiations) and the four most offensive players (scoring opportunities) as a second group (in the sense of Definition 2) into the analysis. The clusterings in the remainder are thus based on the trajectories of five objects (ball, four players).

In general, there is no ground-truth for real-world clustering problems. Figure 5 therefore depicts the adjusted Rand indices between the five methods for pair-wise comparison. Generally, the adjusted Rand index is low in the majority of the cases indicating low consistency between the methods. On the other hand, we have a higher consistency between STCK<sub>dist</sub> and DTW<sub>dist</sub> demonstrating that our method is capable of dealing with speed and length differences in a similar way as dynamic time warping. This is also partially due to the use of the same kernel in both methods, once as spatial kernel and once as local distance measure.

Figure 6 (left) shows average Silhouette measures of the four similaritybased methods over the four datasets. Our STCK<sub>dist</sub> achieves the highest score indicating higher cluster separation and/or compactness. The object-wise STCK provides the poorest results on average, indicating that permutations are relevant in this application and that the distribution-based representation of the probability product kernel is more successful in capturing the relevant player movements. Dynamic time warping performs second best on average and in line with the previous results on the artificial datasets. Also note the high consistency with our method in Figure 5. In terms of 5-nearest-neighbour consistency, the two distribution-based methods,  $STCK_{dist}$  and  $DTW_{dist}$ , perform best, as depicted in Figure 6 (right). Cluster quality is generally higher for game initiations compared to scoring opportunities which also results in in more interpretable clusters for these settings as we will show in the remainder.

Figures 7 and 8 illustrate the resulting clusters for the five methods. The medoids of the clusters for both teams are depicted along with the distribution of the trajectory length for each cluster. For all settings cluster membership correlates with sequence length. Generally, the cluster medoids of  $\text{STCK}_{\text{dist}}$  and  $\text{DTW}_{\text{dist}}$  often coincide and are easily interpretable, while the cluster medoids of  $\text{STCK}_{1-1}$  and the cluster representatives of the multinomial mixture model are difficult to make sense of.



Only the trajectory of the ball is shown, while the clusters are also based on the trajectories of four players. Where possible, similar clusters are coloured equally. Below, the distribution of the length of the trajectories in each cluster is depicted. For multinomial mixture models, the multi-object trajectory with the highest probability is depicted as representative. Fig. 7 Game initiation clusters of Team A (top; 4 Clusters) and Team B (bottom; 6). The above plots show cluster medoids indicated by colours.



Fig. 8 Scoring opportunity clusters of Team A (top; 6 Clusters) and Team B (bottom; 10). The plots show the cluster medoids. Colours indicate clusters. Only the trajectory of the ball is shown, while the clusters are also based on the trajectories of four players. Where possible, similar clusters are coloured equally. Below, the distribution of the length of the trajectories in each cluster is depicted. In the case of the multinomial mixture model, the multi-object trajectory with the highest probability is depicted as representative.

In the Bundesliga 2011/2012 season, the strategy of Team A generally consisted of transporting the ball with few, but rehearsed short game initiations to the opposing half. For this purpose, many ball contacts were allowed and different players were integrated. On the contrary, Team B showed a rather chaotic game organisation with rather random actions and increasingly long, straight balls. Figure 7 shows the outcomes of the k-medoids cluster for the game initiations. Compared to all other methods, the the STCK<sub>dist</sub> capture the characteristic traits of the teams well. Team A clearly acted with many short moves (long trajectories in cluster 1) and integrated many players in the playmaking (cluttered medians). By contrast, Team B acted with many long moves (short trajectories in all clusters) and preferred linear actions.

Near the opposing goal, Team A aimed at quickly achieving a goal in the opposing half during the 2011/2012 season. They operated with only a few ball contacts and aimed to quickly transport the ball in the predefined zone of danger. Again in contrast to this, Team B had many ball contacts and took their time in waiting for a mistake of the opponent and only then played in the zone of danger to achieve a goal. Figure 8 shows that all methods are capable of retracing the different offensive strategies of both teams. Again, Team B has more solution categories (more than 33%) than Team A, which is shown by the versatile and multifaceted running patterns. However, solely the results with STCK<sub>dist</sub> show that Team A rapidly tries to enter the zone of danger to Team B). To sum up, the results with STCK<sub>dist</sub> best reflect the game philosophies of Team A and B from a sport-scientific perspectives and are the most easy to interpret.

#### **5** Conclusion

We presented spatio-temporal convolution kernels for multi-object scenarios. Our kernels consist of a temporal and a spatial component that can be chosen according to the characteristic traits of a problem at-hand. The computation time is quadratic in terms of the number and lengths of trajectories. We proposed an efficient percental approximation algorithm that significantly reduced the complexity to superlinear runtime. Empirical results on artificial clustering tasks showed that our spatio-temporal convolution kernels effectively identify the target concepts. Results on large-scale real world data from soccer games showed that our kernels lead to easily interpretable clusters that may be used in further analysis by coaches.

#### References

- 1. Andrieu, C.: An introduction to MCMC for machine learning (2003)
- Basharat, A., Gritai, A., Shah, M.: Learning object motion patterns for anomaly detection and improved object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (2008)

- Baud, O., El-Bied, Y., Honore, N., Taupin, O.: Trajectory comparison for civil aircraft. In: Aerospace Conference, 2007 IEEE (2007)
- Bellman, R., Kalaba, R.: On adaptive control processes. IRE Transactions on Automatic Control 4(2) (1959)
- 5. Białkowski, A., Lucey, P., Carr, P., Denman, S., Matthews, I., Sridharan, S.: Recognising team activities from noisy data. In: IEEE Conference on Computer Vision and Pattern Recognition Workshops (2013)
- 6. Białkowski, A., Lucey, P., Carr, P., Yue, Y., Matthews, I.: "Win at Home and Draw Away": Automatic Formation Analysis Highlighting the Differences in Home and Away Team Behaviors. In: MIT Sloan Sports Analytics Conference (2014)
- Blanchard, G., Bousquet, O., Massart, P.: Statistical performance of support vector machines. The Annals of Statistics 36(2) (2008)
- Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. Journal for Machine Learning Research 3 (2003)
- 9. Buzan, D., Sclaroff, S., Kollios, G.: Extraction and clustering of motion trajectories in video. In: International Conference on Pattern Recognition, vol. 2 (2004)
- Chen, Y., Wiesel, A., Eldar, Y., Hero, A.: Shrinkage algorithms for MMSE covariance estimation. IEEE Transactions on Signal Processing 58(10) (2010)
- Direkoglu, C., O'Connor, N.: Team behavior analysis in sports using the poisson equation. In: IEEE International Conference on Image Processing (2012)
- Fowlkes, C., Belongie, S., Chung, F., Malik, J.: Spectral grouping using the Nyström method. IEEE Transactions on Pattern Analysis and Machine Intelligence 26(2) (2004)
- Fu, Z., Hu, W., Tan, T.: Similarity based vehicle trajectory clustering and anomaly detection. In: International Conference on Image Processing, vol. 2 (2005)
- Giannotti, F., Nanni, M., Pinelli, F., Pedreschi, D.: Trajectory pattern mining. In: Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '07. ACM, New York, NY, USA (2007). DOI 10.1145/1281192.1281230
- 15. Grunz, A., Memmert, D., Perl, J.: Tactical pattern recognition in soccer games by means of special self-organizing maps. Human movement science (2012)
- Hartigan, J.A.: Clustering Algorithms, 99th edn. John Wiley & Sons, Inc., New York, NY, USA (1975)
- 17. Hausdorff, F.: Set theory. Chelsea Pub. Co. (1962)
- Haussler, D.: Convolution kernels on discrete structures. Technical Report UCS-CRL-99-10, University of California at Santa Cruz (1999)
- Hervieu, A., Bouthemy, P.: Understanding sports video using players trajectories. In: Intelligent Video Event Analysis and Understanding, *Studies in Computational Intelli*gence, vol. 332. Springer Berlin Heidelberg (2010)
- Hirano, S., Tsumoto, S.: Grouping of soccer game records by multiscale comparison technique and rough clustering. In: International Conference on Hybrid Intelligent Systems (2005)
- Hu, W., Xiao, X., Fu, Z., Xie, D., Tan, T., Maybank, S.: A system for learning statistical motion patterns. IEEE Transactions on Pattern Analysis and Machine Intelligence 28(9) (2006)
- 22. Hubert, L., Arabie, P.: Comparing partitions. Journal of Classif. 2(1), 193-218 (1985)
- Impire AG: VIS.TRACK Produktinformation (2014). URL http://www.bundesligadatenbank.de/fileadmin/impire/products/Produktinformationen/VIS.TRACK\_-Produktionformation.pdf. [Online; accessed 10-June-2014]
- Intille, S.S., Bobick, A.F.: Recognizing planned multiperson action. Computer Vision and Image Understanding 81(3) (2001)
- Jebara, T., Kondor, R., Howard, A.: Probability product kernels. Journal for Machine Learning Research 5 (2004)
- Jeong, H., Chang, H.J., Choi, J.Y.: Modeling of moving object trajectory by spatiotemporal learning for abnormal behavior detection. In: IEEE International Conference on Advanced Video and Signal-Based Surveillance (2011)
- Jinyang, D., Rangding, W., Liangxu, L., Jiatao, S.: Clustering of trajectories based on hausdorff distance. In: International Conference on Electronics, Communications and Control (2011)

- Junejo, I., Javed, O., Shah, M.: Multi feature path modeling for video surveillance. In: International Conference on Pattern Recognition, vol. 2 (2004)
- Kang, C.H., Hwang, J.R., Li, K.J.: Trajectory analysis for soccer players. In: IEEE International Conference on Data Mining Workshops (2006)
- Kang, J., Yong, H.S.: Spatio-temporal discretization for sequential pattern mining. In: Proceedings of the 2Nd International Conference on Ubiquitous Information Management and Communication. ACM, New York, NY, USA (2008)
- Kang, S.J., Kim, Y., Park, T., Kim, C.H.: Automatic player behavior analysis system using trajectory data in a massive multiplayer online game. Multimedia Tools and Applications (2013). DOI 10.1007/s11042-012-1052-x
- Kaufman, L., Rousseeuw, P.: Clustering by means of medoids. Statistical Data Analysis Based on the L1-Norm and Related Methods (1987)
- Kempe, M., Grunz, A., Memmert, D.: Detecting tactical patterns in basketball: Comparison of merge self-organising maps and dynamic controlled neural networks. European Journal of Sport Science (2014)
- Kim, K., Grundmann, M., Shamir, A., Matthews, I., Hodgins, J., Essa, I.: Motion fields to predict play evolution in dynamic sport scenes. In: IEEE Conference on Computer Vision and Pattern Recognition (2010)
- Larson, J.S., Bradlow, E.T., Fader, P.S.: An exploratory look at supermarket shopping paths. International Journal of Research in Marketing 22 (2005). DOI 10.1016/j.ijresmar.2005.09.005
- 36. Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In: IEEE Computer Society Conference on Computer Vision and Pattern Recognition, vol. 2 (2006)
- Ledoit, O., Wolf, M.: A well-conditioned estimator for large-dimensional covariance matrices. Journal of Multivariate Analysis 88(2) (2004)
- Lee, J.G., Han, J., Li, X.: Trajectory outlier detection: A partition-and-detect framework. In: Proceedings of the 2008 IEEE 24th International Conference on Data Engineering. IEEE Computer Society, Washington, DC, USA (2008). DOI 10.1109/ICDE.2008.4497422
- Li, R., Chellappa, R.: Group motion segmentation using a spatio-temporal driving force model. In: IEEE Conference on Computer Vision and Pattern Recognition (2010)
- Li, R., Chellappa, R., Zhou, S.: Learning multi-modal densities on discriminative temporal interaction manifold for group activity recognition. In: IEEE Conference on Computer Vision and Pattern Recognition (2009)
- Lin, C.J.: On the convergence of the decomposition method for support vector machines. IEEE Transactions on Neural Networks 12(6) (2001)
- 42. Lin, D., Grimson, E., Fisher, J.: Learning visual flows: A lie algebraic approach. In: IEEE Conference on Computer Vision and Pattern Recognition (2009)
- Lucey, P.J., Bialkowski, A., Carr, P., Morgan, S., Matthews, I., Sheikh, Y.: Representing and discovering adversarial team behaviors using player roles. In: IEEE Conference on Computer Vision and Pattern Recognition (2013)
- 44. Mamoulis, N., Cao, H., Kollios, G., Hadjieleftheriou, M., Tao, Y., Cheung, D.W.: Mining, indexing, and querying historical spatiotemporal data. In: Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, New York, NY, USA (2004)
- 45. Mehta, S., Parthasarathy, S., Yang, H.: Toward unsupervised correlation preserving discretization. Knowledge and Data Engineering, IEEE Transactions on (2005)
- 46. Ng, A.Y., Jordan, M.I., Weiss, Y.: On spectral clustering: Analysis and an algorithm. In: Advances in Neural Information Processing Systems. MIT Press (2001)
- 47. Ong, C.S., Canu, S., Smola, A.J.: Learning with non-positive kernels. In: International Conference on Machine Learning (2004)
- Pao, H.K., Chen, K.T., Chang, H.C.: Game bot detection via avatar trajectory analysis. Computational Intelligence and AI in Games, IEEE Transactions on (2010)
- Perše, M., Kristan, M., Kovačič, S., Vučkovič, G., Perš, J.: A trajectory-based analysis of coordinated team activity in a basketball game. Computer Vision and Image Understanding 113(5) (2009)

- Piciarelli, C., Foresti, G., Snidaro, L.: Trajectory clustering and its applications for video surveillance. In: IEEE Conference on Advanced Video and Signal Based Surveillance (2005)
- 51. Rand, W.M.: Objective criteria for the evaluation of clustering methods. Journal of the American Statistical Association 66(336) (1971)
- 52. Rousseeuw, P.J.: Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. Journal of Computational and Applied Mathematics **20** (1987)
- Saleemi, I., Shafique, K., Shah, M.: Probabilistic modeling of scene dynamics for applications in visual surveillance. IEEE Transactions on Pattern Analysis and Machine Intelligence 31(8) (2009)
- Shawe-Taylor, J., Cristianini, N.: Kernel Methods for Pattern Analysis. Cambridge University Press (2004)
- 55. Siddiquie, B., Yacoob, Y., Davis, L.: Recognizing plays in american football videos
- Steinwart, I.: Support vector machines are universally consistent. Journal of Complexity 18(3) (2002)
- Trunfio, G.A., D'Ambrosio, D., Rongo, R., Spataro, W., Di Gregorio, S.: A new algorithm for simulating wildfire spread through cellular automata. ACM Trans. Model. Comput. Simul. (2011)
- Vapnik, V.N.: The Nature of Statistical Learning Theory. Springer-Verlag New York, Inc., USA (1995)
- Wang, X., Ma, K.T., Ng, G.W., Grimson, W.: Trajectory analysis and semantic region modeling using a nonparametric bayesian model. In: IEEE Conference on Computer Vision and Pattern Recognition (2008)
- Wauthier, F.L., Jojic, N., Jordan, M.I.: Active spectral clustering via iterative uncertainty reduction. In: ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM (2012)
- Wei, X., Sha, L., Lucey, P., Morgan, S., Sridharan, S.: Large-scale analysis of formations in soccer. In: International Conference on Digital Image Computing: Techniques and Applications (2013)
- Williams, C., Seeger, M.: Using the Nyström method to speed up kernel machines. In: Advances in Neural Information Processing Systems 13. MIT Press (2001)
- Zhu, G., Huang, Q., Xu, C., Rui, Y., Jiang, S., Gao, W., Yao, H.: Trajectory based event tactics analysis in broadcast sports video. In: International Conference on Multimedia. ACM (2007)

## A Generation of Artificial Data

The four clusters of the linear movement data consist of trajectories each of which has of five objects performing a linear movement as depicted in Figure 1. The clusters differ with regard to the direction of the movement as follows:

- Cluster 1: Objects move downwards in parallel. In particular, the position (without noise) of object  $o_i$  for i = 1, ..., 5 at time t = 0, 0.1, ..., 1 is given by

$$\left(-0.2 + \frac{i-1}{10}, 1-t\right).$$

– Cluster 2: Objects move upwards in parallel. In particular, the position of object  $o_i$  for i = 1, ..., 5 at time t = 0, 0.1, ..., 1 is given by

$$\left(-0.2+\frac{i-1}{10},t\right).$$

– Cluster 3: Objects move downwards while drifting apart. In particular, the position of object $o_i$  for i=1,2 at time t=0,0.1,...,1 is

$$(-0.6 + (i-1) \cdot 0.1 + 0.4 \cdot t, 1 - t),$$

the position of object  $o_3$  is

$$(0, 1-t)$$

at time t = 0, 0.1, ..., 1 and the position of object  $o_i$  for i = 4, 5 at time t = 0, 0.1, ..., 1 is

$$(0.5 + (i - 4) \cdot 0.1 - 0.4 \cdot t, 1 - t).$$

 – Cluster 4: Objects move upwards while drifting apart. In particular, the position of object<br/>  $o_i$  for i=1,2 at time t=0,0.1,...,1 is

$$(-0.6 + (i-1) \cdot 0.1 + 0.4 \cdot (1-t), t),$$

the position of object  $o_3$  is

(0,t) at time t=0,0.1,...,1 and the position of object  $o_i$  for i=4,5 at time t=0,0.1,...,1 is

$$(0.5 + (i - 4) \cdot 0.1 - 0.4 \cdot (1 - t), t).$$

The second artificial data set testing on circular movements contains trajectories of four objects moving in circles of different radii. The clusters differ with regard to the orientation of the movement as follows:

- Cluster 1: All objects move clockwise. In particular, the position of object  $o_i$  for i = 1, ..., 4 at time t = 0, 0.1, ..., 1 is given by

$$(0.25i\cos(2\pi t), 0.25i\sin(2\pi t)).$$

– Cluster 2: All objects move counter-clockwise. In particular, the position of object  $o_i$  for i = 1, ..., 4 at time t = 0, 0.1, ..., 1 is given by

 $(0.25i, \cos(2\pi t), -0.25i\sin(2\pi t)).$ 

- Cluster 3: The inner two objects move clockwise, the outer two counter-clockwise. In particular, the position of object  $o_i$  for i = 1, ..., 4 at time t = 0, 0.1, ..., 1 is given by

 $(0.25i, \cos(2\pi t), \operatorname{sgn}(2-i)0.25i\sin(2\pi t)).$ 

- Cluster 4: The inner two objects move counter-clockwise, the outer two clockwise. In particular, the position of object  $o_i$  for i = 1, ..., 4 at time t = 0, 0.1, ..., 1 is given by

 $(0.25i, \cos(2\pi t), \operatorname{sgn}(i-3)0.25i\sin(2\pi t)).$