

Mining Positional Data Streams

Jens Haase and Ulf Brefeld

Knowledge Mining & Assessment Group
TU Darmstadt and DIPF
brefeld@cs.tu-darmstadt.de

Abstract. We study frequent pattern mining from positional data streams. Existing approaches require discretised data to identify atomic events and are not applicable in our continuous setting. We propose an efficient trajectory-based preprocessing to identify similar movements and a distributed pattern mining algorithm to identify frequent trajectories. We empirically evaluate all parts of the processing pipeline.

1 Introduction

Recent advances in telecommunication, sensing, and recording technologies allow for storing positions from moving objects at large scales in (near) real time. Analysing positional data streams is highly important in many applications; examples range from navigation and routing systems, network traffic, animal migration/tracking, movements of avatars in computer games to tactics in team sports.

In this paper, we focus on identifying frequent movement patterns in positional data streams that consist of a possible infinite sequence of coordinates. Existing approaches to frequent pattern mining [3, 29, 20] use identities of atomic events to define sequences (episodes) [17, 15]. In positional data, events correspond to sequences of positions (i.e., trajectories) and due to the continuous domain it is very unlikely to observe a trajectory twice. Instead, we observe a multitude of different trajectories that give rise to an exponentially growing set of possibly frequent sequences. Consequently, mining positional data can only be addressed in the context of big data.

Our contribution is threefold: (i) To remedy the absence of matching atomic events, we propose an efficient preprocessing of the positional data using locality sensitive hashing and approximate dynamic time warping. (ii) To process the resulting near-neighbour trajectories we present a novel frequent pattern mining algorithm that generalises Achar et al. [1] to positional data. (iii) We present a distributed algorithm for processing positional data at large-scales. Empirically, we evaluate all stages of our approach on positional data of a real soccer game where cameras and sensors realise a bird's eye view of the pitch that allows for locating the players and the ball several times per second.

The next section reviews related work. Section 3 introduces the representation and the efficient computation of near neighbours. We present our algorithm to detect frequent trajectories in Section 4. Section 5 reports on empirical results and Section 6 concludes.

2 Related Work

Spatio-temporal data mining aims to extract the behaviour and relation of moving objects from (positional) data streams and is frequently used in computational biology for mining animal movements. In [11], the authors aim to detect closely moving objects. Although composition and location of the group may change over time, the identity of the group is considered fixed. Such groups of objects are called *moving clusters*.

Trajectory-based patterns are first introduced by [6]. These patterns represent a set of individual trajectories that share the property of visiting the same sequence of places within similar travel time. Trajectory-based approaches use a discretisation of the movements to identify places that are also known beforehand. Our contribution considers a continuous generalisation: every coordinate on the pitch is a place of interest and trajectories are relations between coordinates and travel time.

Event sequence mining has been introduced by [3] as a problem of mining frequent sequential patterns in a set of sequences. Sequential pattern mining discovers frequent subsequences as patterns in a sequence database. The most common example is the cart analysis proposed by [3]. Their approach discovers items that are often bought together by customers. Efficient generalisations have been proposed to mine large sequence databases [29, 20]. A special case of sequential pattern mining that aims at closed patterns is introduced in [28]. Frequent episode discovery algorithms [17] summarise techniques to describe and find patterns in a stream of events. Applications in several domains have been developed like manufacturing [15], telecommunication [17], biology [19], and text mining [9]. However, these algorithms are restricted to specific types of episodes. Achar et al. [1] propose the first approach to mine unrestricted episodes. Our approach generalises [1] to mining positional data streams.

Data mining for team sports mostly focuses on event recognition in videos [5, 27] or the conversion of video data into a positional data stream [10]. [10] also report on simple descriptive statistics including heat maps and passing accuracies. [12] uses positional data to assess player positions in particular areas of the pitch, such as catchable, safe or competing zones. Prior work also utilises positional data of basketball games [21] and military operations [24] to identify tactical patterns, respectively. However, the presented approaches focus on detecting *a priori* known patterns in the data stream. By contrast, we devise a purely data-driven approach to find frequent patterns in positional data without making any assumptions on zones, tasks or movements.

3 Efficiently Finding Similar Movements

3.1 Representation

Given a positional data stream \mathcal{D} with ℓ objects $\mathbf{o}_1, \dots, \mathbf{o}_\ell$. Every object \mathbf{o}_i is represented by a sequence of coordinates $\mathcal{P}_i = \langle \mathbf{x}_1^i, \mathbf{x}_2^i, \dots \rangle$ where $\mathbf{x}_t = (x_1, x_2, \dots, x_d)^\top$ denotes the position of the object in d -dimensional space at

time t . A trajectory or movement of the i -th object is a subset $\mathbf{p}_{[t,t+m]} \subseteq \mathcal{P}_i$ of the stream, e.g., $\mathbf{p}_{[t,t+m]} = \langle \mathbf{x}_t^i, \mathbf{x}_{t+1}^i, \dots, \mathbf{x}_{t+m}^i \rangle$, where m is the length of the trajectory. In the remainder, the time index t is omitted and each element of a trajectory is indexed by offsets $1, \dots, m$.

For generality, we focus on finding similar trajectories where (i) the exact location of a trajectory does not matter (*translation invariance*), (ii) the range of the trajectory is negligible (*scale invariance*), and where turns such as left or right are considered identical (*rotation invariance*). Note that, depending on the application at hand, one or more of these requirements may be inappropriate and can be dropped by altering the representation accordingly. Using the requirements (i)-(iii) gives rise to the so-called angle/arc-length representation [26] of trajectories that represents movements as a list of tuples of angles θ_t and distances $\mathbf{v}_t = \mathbf{x}_t - \mathbf{x}_{t-1}$. The difference \mathbf{v}_t is called the *movement vector* at time t and the angles are computed with respect to a (randomly drawn) reference vector $\mathbf{v}_{ref} = (1, 0)^\top$,

$$\theta_i = \text{sign}(\mathbf{v}_i, \mathbf{v}_{ref}) \left[\cos^{-1} \left(\frac{\mathbf{v}_i^\top \mathbf{v}_{ref}}{\|\mathbf{v}_i\| \|\mathbf{v}_{ref}\|} \right) \right],$$

where the sign function computes the direction (clockwise or counterclockwise) of the movement with respect to the reference.

Additionally, transformed trajectories are normalised by subtracting the average so that $\theta_i \in [-\pi, +\pi]$ for all i and by normalising the total distance to one. Finally, we discard the difference vectors and represent trajectories solely by their sequences of angles, $\mathbf{p} \mapsto \tilde{\mathbf{p}} = \langle \theta_1, \dots, \theta_n \rangle$.

3.2 Approximate Dynamic Time Warping

Recall that pairs of trajectories may contain phase shifts, that is, a movement may begin slowly and then speeds-up while another starts fast and then slows down towards the end. Such phase shifts are well captured by alignment-based similarity measures such as dynamic time warping [22].

Dynamic time warping (DTW) is a non-metric distance function that measures the distance between two sequences and is often used to compare time related signals in biometrics or speech recognition problems. Given two sequences $\mathbf{s} = \langle s_1, \dots, s_n \rangle$ and $\mathbf{q} = \langle q_1, \dots, q_m \rangle$ and a cost function $\text{cost}(s_i, q_j)$ detailing the costs of matching s_i with q_j . The goal of dynamic time warping is to find an alignment $A = \{(\eta_i, \mu_i)\}$ for $\eta_i \in [0, n]$ and $\mu_i \in [0, m]$ of sequences \mathbf{s} and \mathbf{q} that has minimal costs subject to the constraints that (i) the alignment starts at position $(1, 1)$ and ends at position (n, m) (*boundary*), (ii) being in step (η_i, μ_j) after being in step (η_k, μ_l) implies that $i - k \leq 1$ and $j - l \leq 1$ (*continuity*), and (iii) being in step (η_i, μ_j) after being in step (η_k, μ_l) implies that $i - k \geq 0$ and $j - l \geq 0$ (*monotonicity*). An alignment that fulfils the above conditions is given by [18] using $g(\emptyset, \emptyset) = 0$, $g(\mathbf{s}, \emptyset) = \text{cost}(\emptyset, \mathbf{q}) = \infty$, and

$$g(\mathbf{s}, \mathbf{q}) = \text{cost}(s_1, q_1) + \min \left\{ \begin{array}{l} g(\mathbf{s}, \langle q_2, \dots, q_m \rangle) \\ g(\langle s_2, \dots, s_m \rangle, \mathbf{q}) \\ g(\langle s_2, \dots, s_m \rangle, \langle q_2, \dots, q_m \rangle) \end{array} \right\}.$$

The cost function *cost* is arbitrary and any metric or non-metric distance functions can be used. Dynamic time warping has a complexity of $\mathcal{O}(|\mathbf{s}||\mathbf{q}|)$ which is prohibitive for mining positional data streams.

A great deal of methods has been proposed to speed-up the computation of dynamic time warping. Besides global constraints (e.g., [23, 8]), efficient approximations can be obtained by lower bounds. The rationale is that lower bound functions can be computed in less time and are therefore often used as pruning techniques in applications like indexing or information retrieval. The exact DTW computation only needs to be carried out if the lower bound value is above a given threshold. We make use of two lower bound functions, f_{kim} [14] and f_{keogh} [13], that are defined as follows: f_{kim} focuses on the first, last, greatest and smallest values of two sequences [14] and can be computed in $\mathcal{O}(m)$:

$$f_{kim}(\mathbf{s}, \mathbf{q}) = \max \{ |s_1 - q_1|, |s_m - q_m|, |\max(\mathbf{s}) - \max(\mathbf{q})|, |\min(\mathbf{s}) - \min(\mathbf{q})| \}.$$

If the greatest and smallest entries are normalised to a specific value their computation can be ignored and the time complexity reduces to $\mathcal{O}(1)$. The second lower bound f_{keogh} [13] uses minimum ℓ_i and maximum values u_i for sub-sequences of the query \mathbf{q} given by

$$\ell_i = \min(q_{i-r}, \dots, q_{i+r}) \quad \text{and} \quad u_i = \max(q_{i-r}, q_{i+r}),$$

where r is a user defined threshold. Trivially, $u_i \geq q_i \geq \ell_i$ holds for all i and the lower bound f_{keogh} is given by

$$f_{keogh}(\mathbf{q}, \mathbf{s}) = \sqrt{\sum_{i=1}^m c_i} \quad \text{with} \quad c_i = \begin{cases} (s_i - u_i)^2 & : \text{if } s_i > u_i \\ (s_i - \ell_i)^2 & : \text{if } s_i < \ell_i \\ 0 & : \text{otherwise} \end{cases}$$

which can also be computed in $\mathcal{O}(m)$. The result is a non-metric distance function that only violates the triangle inequality of a metric distance.

3.3 An N -Best Algorithm

Given a trajectory $\mathbf{q} \in \mathcal{D}$, the goal is to find the most similar trajectories in \mathcal{D} . Trivially, a straight forward approach is to compute the DTW values of \mathbf{q} for all trajectories in \mathcal{D} and sort the outcomes accordingly. However, this requires $|\mathcal{D}|$ DTW computations, each of which is quadratic in the length of the trajectories, and renders the approach clearly infeasible. Algorithm 1 computes the N most similar trajectories for a given query \mathbf{q} efficiently by making use of the lower bound functions f_{kim} and f_{keogh} . Lines 2–9 compute the DTW distances of the first N entries in the database and stores the entry with the highest distance to \mathbf{q} . Lines 10–21 loop over the other trajectories in \mathcal{D} by first applying the lower bound functions f_{kim} and f_{keogh} to efficiently filter irrelevant movements before using the exact DTW distance for the remaining candidates. Every trajectory, realising a smaller DTW distance than the current maximum, replaces its peer; *maxdist* and *maxind* are updated accordingly. Note that the

Algorithm 1 N -Best($N, \mathbf{q}, \mathcal{D}$)**Input:** number of near-neighbours N , query trajectory \mathbf{q} , stream \mathcal{D} **Output:** The N most similar trajectories to \mathbf{q} in \mathcal{D}

```

1:  $output = \emptyset; maxdist = 0; maxind = -1$ 
2: for  $i = 1, \dots, N$  do
3:    $dist = DTW(\mathbf{q}, \mathcal{D}[i])$ 
4:    $output[i] = \mathcal{D}[i]$ 
5:   if  $dist > maxdist$  then
6:      $maxdist = dist; maxind = i$ 
7:   end if
8: end for
9: for  $i = N + 1, \dots, |\mathcal{D}|$  do
10:  if  $f_{kim}(\mathbf{q}, \mathcal{D}[i]) < maxdist$  then
11:    if  $f_{keogh}(\mathbf{q}, \mathcal{D}[i]) < maxdist$  then
12:       $dist = DTW(\mathbf{q}, \mathcal{D}[i])$ 
13:      if  $dist < maxdist$  then
14:         $output[maxind] = \mathcal{D}[i]$ 
15:         $maxdist = \max\{DTW(\mathbf{q}, output[j]) : 1 \leq j \leq N\}$ 
16:         $maxind = \operatorname{argmax}_j DTW(\mathbf{q}, output[j])$ 
17:      end if
18:    end if
19:  end if
20: end for

```

complexity of Algorithm 1 is linear in the number of trajectories in \mathcal{D} . In the worst case, the sequences are sorted in descending order by the DTW distance, which requires to compute all DTW distances. In practice, however, much lower run-times are observed. A crucial factor is the tightness of the lower bound functions. The better the approximation of the DTW, the better the pruning. For $N = 1$, the maximum value drops faster towards the lowest possible value. By contrast, setting $N = |\mathcal{D}|$ requires to compute the exact DTW distances for all entries in the database. Hence, in most cases, $N \ll |\mathcal{D}|$ is required to reduce the overall computation time. The computation can trivially be distributed with Hadoop; computing distances is performed in the mapper and sorting is done in the reducer.

3.4 Distance-based Hashing

An alternative to the introduced N -Best algorithm provides locality sensitive hashing (LSH) [7]. A general class of LSH functions are called distance-based hashing (DBH) that can be used together with arbitrary spaces and (possibly non-metric) distances [4]. The hash family is constructed as follows. Let $h : \mathcal{X} \rightarrow \mathbb{R}$ be a function that maps elements $x \in \mathcal{X}$ to a real number. Choosing two randomly drawn members $x_1, x_2 \in \mathcal{X}$, the function h is defined as

$$h_{x_1, x_2}(x) = \frac{dist(x, x_1)^2 + dist(x_1, x_2)^2 - dist(x, x_2)^2}{2 dist(x_1, x_2)}.$$

Algorithm 2 FSATransition($\alpha, fsa, t, events$)

```

1: if  $fsa.currentState.Open = \emptyset$  then
2:   return  $fsa$  {FSA is in final state}
3: end if
4: for  $n \in sourceNodes(fsa.currentState.Open)$  do
5:   for  $e \in events$  do
6:     if  $e \sim nodeMapping_\alpha(n)$  then
7:        $fsa.currentState.Open = fsa.currentState.Open \setminus n$ 
8:        $fsa.currentState.Done = fsa.currentState.Done \cup n$ 
9:        $fsa.lastTransition = t$ 
10:      if  $fsa.startTime == undefined$  then
11:         $fsa.startTime = t$ 
12:      end if
13:      break inner loop {Only one possible similarity (injective episode)}
14:    end if
15:  end for
16: end for
17: return  $fsa$ 

```

The binary hash value for x simply verifies whether $h(x)$ lies in an interval $[t_1, t_2]$,

$$h_{x_1, x_2}^{[t_1, t_2]}(x) = \begin{cases} 1 & : h_{x_1, x_2}(x) \in [t_1, t_2] \\ 0 & : otherwise \end{cases},$$

where the boundaries t_1 and t_2 are chosen so that the probability that a randomly drawn $x \in \mathcal{X}$ lies with 50% chance within and with 50% chance outside of the interval. Given the set \mathcal{T} of admissible intervals and function h , the DBH family is now defined as

$$\mathcal{H}_{DBH} = \left\{ h_{x_1, x_2}^{[t_1, t_2]} : x_1, x_2 \in \mathbb{R} \wedge [t_1, t_2] \in \mathcal{T}(x_1, x_2) \right\}.$$

Using random draws from \mathcal{H}_{DBH} , new hash families can be constructed using *AND*- and *OR*-concatenation.

We use DBH to further improve the efficiency of the *N*-Best algorithm by removing a great deal of trajectories before processing them with Algorithm 1. Given a query trajectory $\mathbf{q} \in \mathcal{D}$, the retrieval process first identifies candidate objects that are hashed to the same bucket for at least one of the hash functions, and then computes the exact distances of the remaining candidates using the *N*-Best algorithm. As distance measure of the DBH hash family we use the lower bound f_{kim} . The computation is again easily distributed with Hadoop.

4 Frequent Episode Mining for Positional Data

The main difference between frequent episode mining (e.g., [1]) and mining frequent trajectories from positional data streams is the definition of events. In contrast to having a predefined set of atomic events, every trajectory in the

Algorithm 3 Map(id, α)

```

1:  $eventStream = loadEventStreamFormFile()$ 
2:  $frequency = 0; fsas = \{new\ FSA\}$ 
3: for all  $(t, events) \in eventsStream$  do
4:   for all  $fsa \in fsas$  do
5:      $inStartState = inStartState(fsa)$ 
6:      $hasChanged = FSATransition(\alpha, fsa, t, events)$ 
7:     if  $inStartState$  and  $hasChanged$  then
8:        $fsas = fsas \cup new\ FSA$ 
9:     end if
10:    if  $inFinalState(fsa)$  then
11:       $fsas = \{new\ FSA\}$ 
12:       $frequency+ = 1$ 
13:    else
14:       $fsas = RemoveAllOlderFSAsInSameState(fsas)$ 
15:    end if
16:  end for
17: end for
18: if  $frequency \geq userDefinedThreshold$  then
19:    $EMIT(blockstart - id(\alpha), \alpha)$ 
20: end if

```

stream is considered an event. Thus, events may overlap and are very unlikely to occur more than just once. In the absence of a predefined set of atomic events, we use the previously defined approximate distance functions in the mining step.

An event stream is a time-ordered stream of trajectories. Every event is represented by a tuple (A, t) where A is an event and t denotes its timestamp. An *episode* α is a directed acyclic graph, described by a triplet (\mathcal{V}, \leq, m) where \mathcal{V} is a set of nodes, \leq is a partial order on \mathcal{V} (directed edges between the nodes), and $m : \mathcal{V} \rightarrow E$ is a bijective function that maps nodes to events in the event stream. We focus on *transitive closed episodes* [25] in the remainder, that is if node A is ordered before B ($A < B$) there must be a direct edge between A and B , that is, $\forall A, B \in \mathcal{V}$ if $A < B \implies edge(A, B)$.

The partial ordering of nodes upper bounds the number of possible directed acyclic graphs on the event stream. The ordering makes it impossible to include two identical (or similar) events in the same episode. Episodes that do not allow duplicate events are called *injective episodes* [1]. An episode α is called *frequent*, if it occurs often enough in the event stream. The process of counting the episode α consists of finding all episodes that are similar to α . A sub-episode β of an episode α can be created by removing exactly one node n and all its edges from and to n ; e.g., for the episode $A \rightarrow B \rightarrow C$ the sub-episodes are $A \rightarrow B$, $A \rightarrow C$ and $B \rightarrow C$. The sub-episode of a singleton is denoted by the empty set \emptyset .

Generally, frequent episodes can be found by Apriori-like algorithms [2]. The principles of dynamic programming are exploited to combine already frequent episodes to larger ones [17, 16]. We differentiate between alternating *episode generation* and *counting* phases. Every newly generated episodes must be *unique*,

transitive closed, and *injective*. Candidates possessing infrequent sub-episodes are discarded due to the downward closure lemma [1]. We now present counting and episode generation algorithms for processing positional data with Hadoop.

4.1 Counting phase

The *frequency* of an episode is defined as the maximum number of non-overlapping occurrences of the episode in the event stream [16].¹ Non-overlapping episodes can be detected and counted with finite state automata (FSAs), where every FSA is tailored to accept only a particular episode. The idea is as follows. For every episode that needs to be counted, an FSA is created and the event stream is processed by each FSA. If an FSA moves out of the initial state, a new FSA is created for possibly later occurring episodes and once the final state has been reached, the episode counter is incremented and all FSA-instances of the episode are deleted except for the one still remaining in the initial state.

Algorithm 2 shows the FSA transition function that counts an instance of an episode. Whenever the FSA reaches its final state its frequency is incremented. As input, Algorithm 2 gets the *fsa* instance which contains the current state, the last transition time and the first transition time. Additionally, the appropriate episode, the current time stamp and the events starting at that time stamp are passed to the function. First, in case the FSA is already in the final state, the function returns without doing anything (line 1). Algorithm 2 iterates over all source nodes in the current state and all events that had happened at time t (line 4-5). Whenever there is an event e that is similar to the appropriate event of source node n (line 6), the FSA is traversed to the next state. The algorithm also keeps track of the start time and the last transition time to check the expiry time (line 9 and line 11).

The FSA transition function can be defined as a counting algorithm shown in Algorithm 3 in terms of a *map*-function for the Hadoop/MapReduce framework. The function first loads the event stream² (line 1) and initialises an empty FSA for every episode. Next, the event stream and the FSAs are traversed and passed to the FSA transition function. Whenever an FSA leaves the start state a new FSA must be added to the set of FSAs. This ensures that there is exactly one FSA in a start state. In case an FSA reaches its final state, all other FSAs can be removed and the process starts again with only one FSA in start state. In case more than one FSA reaches the final state, Algorithm 3 removes all but the youngest one in final state as this one has higher chances to meet the expiry time constraints. The test for expiry time is not shown in the pseudo code. Instances violating the expiry time do not contribute to the frequency count. Neither do FSAs that associate overlapping events with the same object. Note that the general idea of the counting algorithm is very similar to [1]. However, due to the different notions of an event, many optimisation do not apply in our case.

¹ Two occurrences of an episode are said to be *non-overlapping*, if no event associated with one appears in between the events associated with the other.

² In practice one would read the event stream block wise instead of loading the whole data at once into memory. We chose the latter for ease of presentation.

Algorithm 4 Align(α, β)

Require: $|nodes(\alpha)| = |nodes(\beta)|$

- 1: $f =$ int array of length $|nodes(\alpha)|$
- 2: $used =$ boolean array of length $|nodes(\alpha)|$
- 3: $n = 0$
- 4: **for** $i = 1$ **to** $|nodes(\alpha)|$ **do**
- 5: $event_{i,\alpha} = m(\alpha)[i]$
- 6: $found = \mathbf{false}$
- 7: **for** $j = 1$ **to** $|nodes(\beta)|$ **do**
- 8: $event_{j,\beta} = m(\beta)[j]$
- 9: **if** (\mathbf{not} $used[j]$) **and** $event_{i,\alpha} \sim event_{j,\beta}$ **then**
- 10: $f[i] = j$
- 11: $used[j] = \mathbf{true}$
- 12: $found = \mathbf{true}$
- 13: **end if**
- 14: **end for**
- 15: **if** $found = \mathbf{false}$ **then**
- 16: $f[i] = -1$
- 17: $increment(n)$
- 18: **end if**
- 19: **end for**
- 20: **return** f, n

Following [1] we also employ bidirectional evidence as frequencies alone are necessary but not sufficient for detecting frequent episodes. The entropy-based bidirectional evidence can be integrated in the counting algorithm, see [1] for details. We omit the presentation here for lack of space.

4.2 Generation phase

Algorithm 5 is designed to efficiently find the indices of the differentiated nodes of two episodes α and β . Therefore, it first tries to find the bijective mapping π , that maps each node (and its corresponding event) of episode α to episode β (line 1). In case such a complete mapping can not be found, π returns only the possible mappings and n contains the number of missing nodes in the mapping (see Algorithm 4). Episodes α and β are combinable, if and only if $n = 1$. The remainder of the algorithm finds the missing node indices by accumulating over the existing indices and by subtracting the accumulated result from the sum of all indices. This little trick finds the missing indices in time $O(n)$. The function returns the node indices that differentiate between α and β . To prevent the computation of Algorithm 5 on all pairs of episodes, each episode is associated with its *block start identifier* [1]. The idea is the following. All generated episodes from an episode α share the same sub-episode. This sub-episode is trivially identical to α as it originates from adding a node to α . The generation step thus takes only those episodes into account that possess the same block start identifier.

Algorithm 5 Combine(α, β)

```

1:  $\pi, n = \text{Align}(\alpha, \beta)$ 
2: if  $n \neq 1$  then
3:   return  $-1$ 
4: end if
5:  $sum_\alpha = 0; sum_\beta = 0$ 
6:  $sum = \frac{|\pi| \times (|\pi| - 1)}{2}$ 
7: for  $i = 1$  to  $|\pi|$  do
8:   if  $\pi[i] \geq -1$  then
9:      $sum_\alpha = sum_\alpha + i$ 
10:     $sum_\beta = sum_\beta + \pi[i]$ 
11:   end if
12: end for
13: return  $(sum - sum_\alpha, sum - sum_\beta)$ 

```

Given two combinable episodes α and β and the differentiated nodes a and b (found by Algorithm 5), it is now possible to combine these episodes to up to three new candidates, as described by [1]. The first candidate originates from adding node b to episode α including all its edges from and to b . The second candidate is generated from the first candidate by adding an edge from node a to node b and the third one adds an edge from b to a to the first candidate. In contrast to [1], we do not test whether all sub-episodes of each candidate are frequent as this would require an efficient lookup of all episodes which can be quite complex for positional data. Candidates with infrequent sub-episodes are allowed at this stage of the algorithm as they will be eliminated in the next counting step anyway.

The complete episode generation algorithm is shown in Algorithm 6. As input, a list of frequent episodes ordered by their block start identifier is given. The result of the algorithm is a list of new episodes that are passed on to the counting algorithm. In line 2 and line 4, all episode pairs are processed as long as they share the same block start identifier (line 6). Then, three possible candidates are generated (line 7) and kept in case they are transitive closed (line 9). Before adding it to the result set, the block start identifier of the new episode is updated (line 10). Analogously to the counting phase, domain specific constraints may be added to filter out unwanted episodes (e.g. in terms of expiry time, overlapping trajectories of the same object, etc.).

5 Empirical Evaluation

5.1 Positional Data

For the experiments, we use positional data from the DEBS Grand Challenge 2013³. The data is recorded from a soccer game of two teams with eight players

³ <http://www.orgs.ttu.edu/debs2013>

Algorithm 6 Reduce(*blockstartId*, *xs*)

```

1:  $k = -1$ ;  $result = \emptyset$ 
2: for  $i = 0$  to  $|xs|$  do
3:    $\alpha = xs(i)$ ;  $currentBlockStart = k + 1$ 
4:   for  $j = i + 1$  to  $|xs|$  do
5:      $\beta = xs(j)$ 
6:     if  $\alpha.blockStart == \beta.blockStart$  then
7:        $candidates = Combine(\alpha, \beta)$ 
8:       for  $c \in candidates$  do
9:         if  $transitiveClose(c)$  then
10:           $c.blockStart = currentBlockStart$ 
11:           $result = result \cup c$ 
12:           $k = k + 1$ 
13:        end if
14:      end for
15:    else
16:       $break$ 
17:    end if
18:  end for
19: end for
20:  $EMIT(id, result)$ 

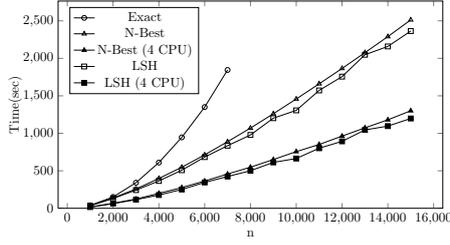
```

on each side. Each player is equipped with two sensors, one for each foot. We average every pair of sensors to obtain only a single measurement for every player at each point in time. Events happening before and after the game or during the half time break are removed as well as coordinates occurring outside of the pitch are discarded. To reduce the overall amount of data, we average the positional data of each player over 100ms. The set of trajectories is created by introducing a sliding window of constant size for each player so that trajectories begin every 500ms and last for one second. This procedure gives us 111.041 trajectories in total, 50.212 for team A, 50.245 for team B, and 10.584 for the ball.

5.2 Near Neighbour Search

The first set of experiments evaluates the run-time of the three distance functions *Exact*, *N-Best*, and *LSH*. Since the exact variant needs quadratically many comparisons in the length of the stream, we focus on only a subset of 15,000 consecutive positions of team A in the experiment. We fix $N = 1000$ and measure the total computation time for all approaches. Figure 1 (left) shows the run-times in seconds for varying sample sizes.

Unsurprisingly, the computation time of the exact distances grows exponentially in the size of the data. By contrast, the *N-Best* algorithm performs slightly super-linear and significantly outperforms its exact counterpart. Pre-filtering trajectories using *LSH* results in only a small additional speed-up. The figure also shows that distributing the computation significantly improves the run-time of the algorithms and indicates that parallelisation allows for computing near-



n	f_{kim}	f_{keogh}	LSH	Σ
1000	0%	0%	11.42%	11.42%
5000	0.28%	34.00%	16.33%	50.61%
10000	9.79%	41.51%	17.80%	60.10%
15000	17.50%	46.25%	11.82%	75.57%

Fig. 1. Left: Run-time. Right: Pruning efficiency.

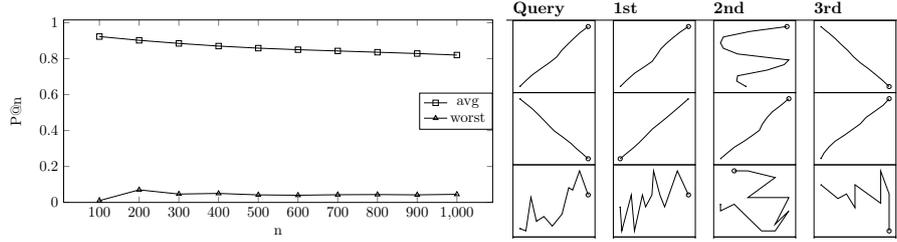


Fig. 2. Left: Accuracy of LSH. Right: Most similar trajectories for a given query

neighbours on large data sets very efficiently. The observed improvements in run-time are the result of a highly efficient pruning strategy. Figure 1 (right) details the amount of pruned trajectories for the respective approaches. The table shows that the effectiveness of the pruning for f_{kim} and f_{keogh} increases with increasing numbers of trajectories. By contrast, pruning LSH is more or less constant and does not change in terms of the number of trajectories but depends on the overall data size and on the ratio $\frac{N}{|D|}$.

We now investigate the accuracy of the proposed LSH pruning. The 1,000 most similar trajectories are compared with the ground truth given by the exact approach. That is, for each trajectory, the two result lists are compared. Figure 2 reports averaged precision@ N scores with $N \in \{100, 200, \dots, 1000\}$ for all trajectories. The average precision decreases slightly for larger N . We credit this observation to the small data set that comprises only a few trajectories. While highly similar trajectories are successfully found, the approximate near-neighbour method fills up remaining slots with non-similar trajectories. The worst precision is therefore close to zero for $N = 100$ and increases slightly for larger N due to only a few true near neighbours in the data.

Although LSH performs well and only slightly decreases in the size of N we focus on the N -Best algorithm with $N = 1,000$ in the remainder for a loss-free and exact computation of the top- N matches. According to Figure 1 (left) the slightly longer execution time is negligible. Figure 2 shows the most similar trajectories for three query trajectories. For common trajectories (top rows), the most similar trajectories are true near neighbours. It can also be seen that the proposed distance function is rotation invariant. For uncommon trajectories (bottom row), the found candidates are very different from the query.

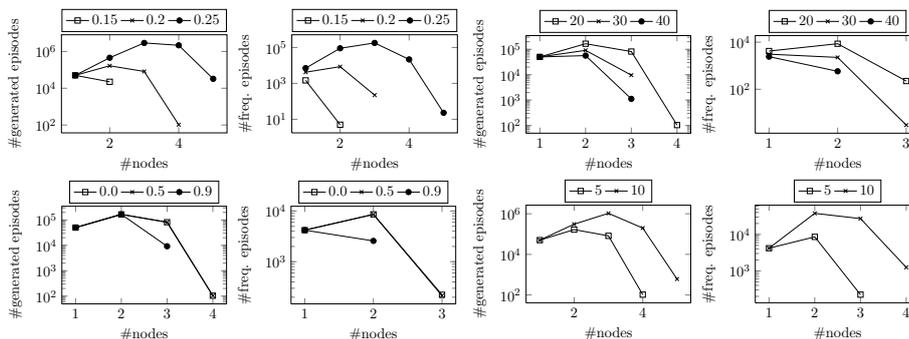


Fig. 3. Top row: Varying similarity (first and second columns) and frequency (third and fourth columns) thresholds. Bottom row: Varying bidirectional evidence (first and second columns) and expiry time (third and fourth columns) thresholds.

5.3 Episode discovery

The first experiments of the episode discovery algorithm focus on the influence of the parameters wrt the number of generated and counted episodes. The algorithm depends on four different parameters, the *similarity*, *frequency*, the *bidirectional evidence*, and the *expiry time*. For this set of experiments, we use the trajectories of team A to find frequent tactical patterns in the game. The results are shown in Figure 3

An interesting parameter is the similarity threshold as it strongly impacts the number of generated episodes: small changes may already lead to an exponential growth in the number of trajectories and large values quickly render the problem infeasible even on medium-sized Hadoop clusters. A similar effect can be observed for the expiry time threshold. Incrementing the expiry time often requires decreasing the similarity threshold. The number of counted episodes is adjusted by the frequency threshold. As shown in the figure, the number of generated episodes can often be reduced by one or more orders of magnitudes. By contrast, the bidirectional evidence threshold affects the result only marginally.

Finally, we present two frequent episodes in Figure 4. The ball is presented with a black line. All other lines describe the players of team A during the time span of the episode. The involved trajectories are displayed by thick lines and a circle at the beginning to indicate movement directions. A small circle at the beginning of a trajectory indicates that the trajectory depends on one or more other trajectories. Additionally, each trajectory is labeled with player ID and timestamp. For completeness, the opponent goal keeper is drawn with a red line at the bottom. The displayed episodes present a well known pattern of soccer games: players move towards the ball. In the left figure, the ball is played towards the opponent goal and players 1, 2, 3, 6, and 7 follow the direction of the ball. In the right figure the opponent team prepares an attack by passing the ball from one side to the other. The players of team A follow the direction of the ball to prevent the attacking team from moving forward.

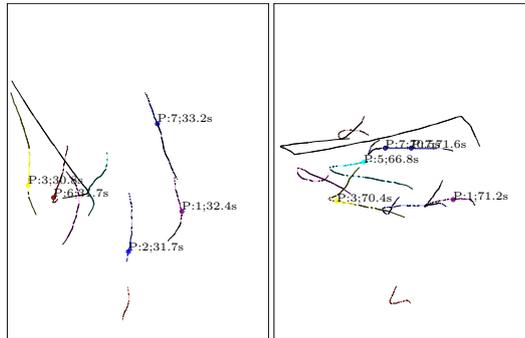


Fig. 4. Exemplary episodes

6 Conclusion

We proposed a novel method to mining frequent patterns in positional data streams where consecutive coordinates of objects are treated as movements. Our contribution is threefold: We firstly proposed an efficient and accurate method to find similar trajectories for a given query. Secondly, we proposed an algorithm that uses these distances to efficiently combine individual movements to complex frequent patterns consisting of multiple trajectories. Thirdly, we presented a distributed version for big data and Hadoop/MapReduce frameworks. Empirical results on positional data from a soccer game showed that the found patterns are intuitive and interpretable.

References

1. A. Achar, S. Laxman, R. Viswanathan, and P. S. Sastry. Discovering injective episodes with general partial orders. *Data Min. Knowl. Discov.*, 25(1):67–108, 2012.
2. R. Agrawal, T. Imieliński, and A. Swami. Mining association rules between sets of items in large databases. *SIGMOD Rec.*, 22(2):207–216, 1993.
3. R. Agrawal and R. Srikant. Mining sequential patterns. In *Proc. of ICDE*, 1995.
4. V. Athitsos, M. Potamias, P. Papapetrou, and G. Kollios. Nearest neighbor retrieval using distance-based hashing. In *Proc. of ICDE*, 2008.
5. M. Beetz, N. v. Hoyningen-Huene, B. Kirchlechner, S. Gedikli, F. Siles, M. Durus, and M. Lames. ASpoGAMo: Automated sports game analysis models. *Int. J. of Computer Science in Sport*, 8(1):4–21, 2009.
6. F. Giannotti, M. Nanni, F. Pinelli, and D. Pedreschi. Trajectory pattern mining. In *Proc. of KDD*, 2007.
7. P. Indyk and R. Motwani. Similarity search in high dimensions via hashing. In *Proc. of VLDB*, 1999.
8. F. Itakura. Minimum prediction residual principle applied to speech recognition. In A. Waibel and K.-F. Lee, editors, *Readings in speech recognition*. Morgan Kaufmann, San Francisco, CA, USA, 1990.

9. K. Iwanuma, Y. Takano, and H. Nabeshima. On anti-monotone frequency measures for extracting sequential patterns from a single very-long data sequence. In *Proc. of CIS*, 2004.
10. S. Iwase and H. Saito. Tracking soccer player using multiple views. In *Proc. of IAPR MVA*, 2002.
11. P. Kalnis, N. Mamoulis, and S. Bakiras. On discovering moving clusters in spatio-temporal data. In *Proc. of SSTD*, 2005.
12. C.-H. Kang, J.-R. Hwang, and K.-J. Li. Trajectory analysis for soccer players. In *Proc. of ICDMW*, 2006.
13. E. Keogh. Exact indexing of dynamic time warping. In *Proc. of VLDB*, 2002.
14. S.-W. Kim, S. Park, and W. W. Chu. An index-based approach for similarity search supporting time warping in large sequence databases. In *Proc. of ICDE*, 2001.
15. S. Laxman. *Discovering frequent episodes: fast algorithms, connections with HMMs and generalizations*. PhD thesis, Indian Institute of Science, 2006.
16. S. Laxman, P. S. Sastry, and K. P. Unnikrishnan. Discovering frequent episodes and learning hidden markov models: A formal connection. *IEEE Trans. on Knowl. and Data Eng.*, 17(11):1505–1517, November 2005.
17. H. Mannila, H. Toivonen, and A. I. Verkamo. Discovery of frequent episodes in event sequences. *Data Min. Knowl. Discov.*, 1(3):259–289, 1997.
18. Meinard Müller. *Information Retrieval for Music and Motion*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2007.
19. D. Patnaik, P. S. Sastry, and K. P. Unnikrishnan. Inferring neuronal network connectivity from spike data: A temporal data mining approach. *Scientific Programming*, 16(1):49–77, 2008.
20. J. Pei, J. Han, B. Mortazavi-asl, H. Pinto, Q. Chen, U. Dayal, and M.C. Hsu. PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth. In *Proc. of ICDE*, 2001.
21. M. Perše, M. Kristan, S. Kovačič, G. Vučkovič, and J. Perš. A trajectory-based analysis of coordinated team activity in a basketball game. *Computer Vision and Image Understanding*, 113(5):612 – 621, 2009.
22. L. Rabiner and B.-H. Juang. *Fundamentals of speech recognition*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1993.
23. H. Sakoe and S. Chiba. Dynamic programming algorithm optimization for spoken word recognition. In A. Waibel and K.-F. Lee, editors, *Readings in speech recognition*, pages 159–165. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1990.
24. G. Sukthankar and K. Sycara. Robust recognition of physical team behaviors using spatio-temporal models. In *Proc. of AAMAS*, 2006.
25. N. Tatti and B. Cule. Mining closed episodes with simultaneous events. In *Proc. of KDD*, 2011.
26. M. Vlachos, D. Gunopulos, and G. Das. Rotation invariant distance measures for trajectories. In *Proc. of KDD*, 2004.
27. C. Xu, Y.-F. Zhang, G. Zhu, Y. Rui, H. Lu, and Q. Huang. Using webcast text for semantic event detection in broadcast sports video. *Trans. Multi.*, 10(7):1342–1355, 2008.
28. X. Yan, J. Han, and R. Afshar. CloSpan: Mining closed sequential patterns in large datasets. In *Proc. of SDM*, 2003.
29. M. J. Zaki. SPADE: An efficient algorithm for mining frequent sequences. *Mach. Learn.*, 42(1-2):31–60, 2001.