The Journal of Systems and Software xxx (2009) xxx-xxx

Contents lists available at ScienceDirect

The Journal of Systems and Software

journal homepage: www.elsevier.com/locate/jss



Searching for similar trajectories in spatial networks

E. Tiakas^a, A.N. Papadopoulos^{a,*}, A. Nanopoulos^a, Y. Manolopoulos^a, Dragan Stojanovic^b, Slobodanka Djordjevic-Kajan^b

^a Department of Informatics, Aristotle University, 54124 Thessaloniki, Greece
^b Department of Computer Science, University of Nis, Aleksandra Medvedeva 14, 18000 Nis, Serbia

ARTICLE INFO

Article history: Received 26 October 2007 Received in revised form 31 October 2008 Accepted 1 November 2008 Available online xxxx

Keywords: Spatial networks Moving objects Trajectories Similarity search

ABSTRACT

In several applications, data objects move on pre-defined spatial networks such as road segments, railways, and invisible air routes. Many of these objects exhibit similarity with respect to their traversed paths, and therefore two objects can be correlated based on their motion similarity. Useful information can be retrieved from these correlations and this knowledge can be used to define similarity classes. In this paper, we study similarity search for moving object trajectories in spatial networks. The problem poses some important challenges, since it is quite different from the case where objects are allowed to move freely in any direction without motion restrictions. New similarity measures should be employed to express similarity between two trajectories that do not necessarily share any common sub-path. We define new similarity in space and time is well expressed, and moreover they satisfy the metric properties. In addition, we demonstrate that similarity range queries in trajectories are efficiently supported by utilizing metric-based access methods, such as M-trees.

© 2008 Elsevier Inc. All rights reserved.

1. Introduction

In location-based services it is important to query the underlying objects based on their location in space, which may change with respect to time. To support such services from the database point of view, specialized tools are required which enable the effective and efficient processing of queries. Queries may involve the spatial or temporal characteristics of the objects, or both (spatio-temporal queries) (Wolfson et al., 1998; Theodoridis et al., 1998). Evidently, indexing schemes are ubiquitous to efficiently support queries on moving objects, by quickly discarding non-relevant parts of the database.

We distinguish between two different research directions towards query processing in moving objects, which differ both in the type of queries supported and the characteristics of the indexing schemes used in each case:

 Query processing techniques for past positions of objects, where past positions of moving objects are archived and queried, using multi-version access methods or specialized access methods for object trajectories (Lomet and Salsberg, 1989; Nascimento and Silva, 1998; Pfoser et al., 2000; Tao and Papadias, 2001a,b). By studying the past positions of objects, important conclusions can be obtained regarding their mobility characteristics. The difficulty in this case is that the database volume increases considerably, since new positions are tracked and recorded.

[II] Query processing techniques for present and future positions of objects, where each moving object is represented as a function of time, giving the ability to determine its future positions according to the current motion characteristics of objects (reference position, velocity vector) (Kollios et al., 1999a,b; Wolfson et al., 2000; Saltenis et al., 2000; Lazaridis et al., 2002). These methods are mainly used to support queries according to the current positions and enable predictions of their future locations. The difficulty in this case is to perform effective predictions, which is difficult taking into consideration that some positions will be invalidated, due to changes in the speed and direction of some objects in the near future.

A data set of moving objects is composed of objects whose positions change with respect to time (e.g., moving vehicles). Since in many cases only the position of each object is important, moving objects are modeled as *moving points* in 2D or 3D Euclidean space. Queries that involve a particular time instance are characterized as *time-slice queries*, whereas queries that must be evaluated for an interval $[t_s, t_e]$ are characterized as *time interval queries*. The

^{*} Corresponding author. Tel.: +30 2310991918; fax: +30 2310991913.

E-mail addresses: tiakas@delab.csd.auth.gr (E. Tiakas), apostol@delab.csd.auth.gr (A.N. Papadopoulos), alex@delab.csd.auth.gr (A. Nanopoulos), manolopo@delab.csd.auth.gr (Y. Manolopoulos), dragans@elfak.ni.ac.yu (D. Stojanovic), sdjordjevic@ elfak.ni.ac.yu (S. Djordjevic-Kajan).

^{0164-1212/\$ -} see front matter \circledcirc 2008 Elsevier Inc. All rights reserved. doi:10.1016/j.jss.2008.11.832

research community has studied both types extensively. Examples of basic queries that could be posed to such a data set include:

- Window query: given a rectangle *R*, which may change position and size with respect to time, determine the objects that are covered by *R* from time point *t*_s to *t*_e.
- *Nearest-neighbor query*: given a moving point *P* determine the *k* nearest-neighbors of *P* within the time interval [*t_s*, *t_e*].
- Join query: given two moving data sets U and V, determine the pairs of objects (o_1, o_2) with $o_1 \in U$ and $o_2 \in V$ such that o_1 and o_2 overlap at some point in $[t_s, t_e]$.

Apart from the query processing techniques proposed for the fundamental types of queries (i.e., window, *k*-NN and join), the issue of *trajectory similarity* has been studied recently. The problem is to identify similar trajectories with respect to a given query trajectory.

The common characteristic of the aforementioned approaches and research works is that objects are allowed to move freely in 2D or 3D space, without any motion restrictions. However, in a large number of applications, objects are allowed to move only on pre-defined paths of an underlying network, resulting in constraint motion. For example, vehicles in a city can only move on road segments. In such a case, the Euclidean distance between two moving objects does not reflect their real distance. Fig. 1 shows such an example which illustrates the differences between restricted and unrestricted trajectories. Objects moving in a spatial network follow specific paths determined by the graph topology, and therefore arbitrary motion is prohibited. This means that two trajectories which are similar regarding the Euclidean distance may be dissimilar when the network distance is considered. The majority of existing methods for trajectory similarity assume that objects can move anywhere in the underlying space, and therefore do not support motion constraints. Most of the proposals are inspired by the time series case, and provide translation invariance, which is not always meaningful in the case of spatial networks. To attack this problem, the network is modeled as a directed graph, and the distance between two objects is evaluated by using algorithms for shortest paths between the nodes of the graph.

Therefore, the challenge is to express trajectory similarity by respecting network constraints, which is also a strong motivation for the following real and practical applications:

- By identifying similar trajectories, effective data mining techniques (e.g., clustering) can be applied to discover useful patterns. For example, a dense cluster is an indication of emerge traffic measures, future road expansions, trafficjam detection, traffic predictions, etc.
- [II] Trajectory similarity can also help in several road network applications such as, routing applications which support historical trajectories, logistic applications, city emergency handling, drive guiding systems, flow analysis, etc. In such applications, efficient indexing and query processing techniques are required.
- [III] Trajectory similarity of moving objects resembles path similarity of user *click-streams* in the area of *web usage mining*. By analyzing the URL path of each user, we are able to determine paths that are very similar, and therefore effective caching strategies can be applied. In web usage mining, web pages and URL links are modeled as a graph. A node in the graph represents a web page, and an edge from one page to another represents an existing link between them. The time spent by each user to a page is also recorded, and it is used in expressing path similarity, in addition to the number of common web pages along each path. In the existing approaches, two paths are considered similar only if they share at least one common web page, or if the paths contain web pages with similar concept. In trajectory similarity on the other hand, two trajectories may be characterized similar even if they do not share any nodes. Therefore, the existing web usage mining techniques are not directly applicable, and the detection of network trajectory similarities can accelerate the web usage mining queries.

The rest of the article is organized as follows. In the next section, we give the appropriate background, we present related work. In Section 3, trajectory similarity search is presented by investigating effective similarity measures between trajectories in a spatial network. Indexing and query processing issues are covered in Section 4, whereas Section 5 offers experimental results. Finally, Section 6 concludes the work.

2. Related work

In several applications, the mobility of objects is constrained by an underlying spatial network. This means that objects cannot





move freely, and their position must satisfy the network constraints. Network connectivity is usually modeled by using a graph representation, composed of a set of vertices (nodes) and a set of edges (connections). Depending on the application, the graph may be *weighted* (a cost is assigned to each edge) and *directed* (each edge has an orientation). Fig. 2 illustrates an example of a spatial network corresponding to a part of a city road network, and its graph representation.

Several research efforts have been performed towards efficient spatial and spatio-temporal query processing in spatial networks. In Sankaranarayanan et al. (2005) nearest-neighbor query processing is achieved by using a mapping technique. This mapping transforms the graph representation of the network to a highdimensional space, where Minkowski metrics can be used. Nearest-neighbor queries in road networks have been also studied in Jensen et al. (2003), where a graph representation is used to model the network. In Papadias et al. (2003) authors study query processing for stationary data sets, by using both a graph representation for the network and a spatial access method. It is shown that the use of Euclidean distance retrieves many candidates, and instead they propose a network expansion method to process range, nearest-neighbor and join queries. In-route nearest-neighbor queries have been studied in Yoo and Shekhar (2005), where given a trajectory source and destination the smallest detour is calculated.

The above contributions deal with efficient spatial or spatiotemporal query processing of fundamental queries like range, nearest-neighbor and join. However, the issue of trajectory similarity has not yet been studied adequately in the case of moving objects in spatial networks. Let T_a and T_b be the trajectories of moving objects o_a and o_b , respectively, and $D(T_a, T_b)$ a function that expresses their dissimilarity in the range [0,1]. If the two objects have similar trajectories we expect the value $D(T_a, T_b)$ to be close to zero. On the other hand, if the two trajectories are completely dissimilar, we expect the value $D(T_a, T_b)$ to be close to one.

An example is illustrated in Fig. 3, where four trajectories are depicted in the 2D Euclidean space. A circle denotes the position of each moving object at the corresponding time instance (t_1, \ldots, t_8) . It is evident that one expects that the two gray-colored trajectories be very similar, in contrast to the two black-colored trajectories.

In several research proposals, trajectory similarity is viewed as the multidimensional counterpart of time series similarity. In Lee et al. (2000), the authors study the problem of similarity search in multidimensional data sequences, to determine similarities in image and video databases. A similarity model based on the Minkowski distance is defined, and each sequence is partitioned to subsequences by means of MBRs, to enable efficient indexing. This work can be viewed as an extension of the method proposed in Faloutsos et al. (1994) for time series data.

In Yanagisawa et al. (2003) a similarity distance between trajectories is defined, which is invariant to translation, rotation and



Fig. 3. Example of four trajectories in the 2D Euclidean space.

scaling. Again, the distance calculation is based on the Minkowski distance. Objects are allowed to move freely in the address space.

In Meratnia and de By (2002), an approach is studied to aggregate similar trajectories using a grid-based spatial unit aggregation. The notion of spatial similarity lies on the neighboring cells of the grid in a standard two-dimensional Euclidean space. Many problems can be arisen with how the grid must be defined, what the cell dimensions must be, and in objects and clusters identification.

In Laurinen et al. (2006), an efficient algorithm for trajectories similarity calculation is presented. But all distance calculations through trajectories are based on Euclidean metrics and spaces (L_p norms).

The method proposed in Vlachos et al. (2002a,b) employs a similarity distance based on the longest common subsequence (LCS) between two trajectories. This approach proposes a distance measure, which is more immune to noise than the Minkowski distance, but does not satisfy the metric space properties, and therefore it is difficult to exploit efficient indexing schemes. Instead, hierarchical clustering is used to group trajectories. Moreover, the similarity measure depends heavily on two parameters, namely δ and ϵ , which must be known in advance, and cannot be altered dynamically without reorganization. These values determine the maximum distance between two locations of different trajectories, in time and space, respectively, to be characterized as similar. Trajectories that differ more are characterized as dissimilar and therefore their similarity is set to zero. This approach does not permit the use of ranking or incremental computation of similarity nearest-neighbor queries.

To the best of the authors' knowledge, the only research work studying trajectory similarity on networks is the work in Hwang et al. (2005, 2006). The authors propose a simple similarity measure based on POIs (points of interest). They retrieve similar trajectories on road network spaces and not in Euclidean spaces. They



Fig. 2. A road network and its graph representation.

propose a filtering method based on spatial similarity and refining similar trajectories based on temporal distance. In order to determine the spatial similarity between trajectories, they define that two trajectories are similar in space by a set of pre-defined points of interest *P* if all points of *P* lie in both trajectories, otherwise they define the two trajectories as dissimilar. There are several drawbacks using this approach:

- The set of points of interest must be pre-defined and controlled by the user which is very restrictive.
- A simple wrong point selection in *P* can harm trajectory spatial similarity and the derived similarity clusters, so points in *P* must be selected very carefully and by an expert of the used road network.
- The similarity in space with such definition (1 = similar, 0 = dissimilar) does not take into account any notion of similarity percentage or similarity range. Therefore, we cannot determine how similar two trajectories are in space.
- The spatial similarity of two trajectories is based only into the fact that they share common points, and not into the general network space. Therefore, many similarities excluded. For example, trajectories that have parallel edges with only a city block distance and no common points, are considered completely dissimilar.

In addition, no details are given with respect to the access methods required to provide efficient similarity search. Moreover, no discussion is performed regarding the metric space properties of the proposed distance measures. Our approach avoids all these drawbacks.

In the sequel, we study in detail the proposed similarity model for trajectory similarity search in spatial networks aiming at: (i) the definition of similarity and distance measures between trajectories that satisfy the metric space properties, (ii) the exploitation of the distance between two graph nodes, which is used as a building block for the definition of trajectory similarity, (iii) the incorporation of time information in the similarity metric, and (iv) the efficient support of similarity queries by exploiting appropriate indexing schemes and applying fast processing algorithms.

3. Trajectory similarity measures

Let \mathcal{T} be a set of trajectories in a spatial network, which is represented by a graph G(V, E), where V is the set of nodes and E the set of edges. Each trajectory $T \in \mathcal{T}$ is defined as:

Table 1

Basic notations used throughout the study.

Symbol	Description		
Τ	Set of trajectories		
8	Set of sub-trajectories		
T, T_a, T_b	Trajectories		
T_q	A query trajectory		
m	Trajectory description length		
G(V, E)	Graph representation of the spatial network		
D_G	Graph diameter		
DE _G	Maximum Euclidean node distance		
v_i	A node in the graph representation		
t _i	Time instance that the object reached node v_i		
е	An edge of the graph		
T[i]. v	The <i>i</i> th node of the trajectory		
T[i].t	The time instance that the object reached the <i>i</i> th node		
$d(v_i, v_j)$	Network-based distance between two nodes		
$de(v_i, v_j)$	Euclidean distance between two nodes		
$D_{netX}(T_a, T_b)$	Network-based distance between trajectories		
$D_{time}(T_a, T_b)$	Time-based distance between trajectories		
Enet	Query radius for network-based similarity		
E _{time}	Query radius for time-based similarity		

$$T = ((v_1, t_1), (v_2, t_2), \dots, (v_m, t_m))$$
(1)

where *m* is the trajectory description length, v_i denotes a node in the graph representation of the spatial network, and t_i is the time instance (expressed in time units, e.g., seconds) that the moving object reached node v_i , and $t_1 < t_i < t_m$, $\forall 1 < i < m$. It is assumed that moving from a node to another comes at a non-zero cost, since at least a small amount of time will be required for the transition. Table 1 gives the most important symbols and the corresponding definitions that are used in our study.

3.1. Expressing trajectory similarity

We will follow a step-by-step construction of the similarity measure by first expressing similarity taking into account only the visited path, ignoring time information. Time information will be considered in a subsequent step.

We begin our exploration by assuming that any two trajectories have the same description length. This assumption will be relaxed later. Let T_a and T_b be two trajectories, each of description length *m*. By using our trajectory definition and ignoring the time information, we have: $T_a = (v_{a1}, v_{a2}, ..., v_{am})$ and $T_b = (v_{b1}, v_{b2}, ..., v_{bm})$, where $\forall i, v_{ai} \in V$ and $v_{bi} \in V$.

Note that, to characterize two trajectories as *similar* it is not necessary that they share common nodes. Therefore, the similarity measure must take into account the *proximity* of the trajectories (how close is one trajectory with respect to the other).

Due to motion restrictions posed by the spatial network, measuring trajectory proximity by means of the Euclidean distance is not appropriate. Instead, it is more natural to use the cost associated with each transition from a graph node to another. For example, in Fig. 4 we observe that two trajectory parts can be similar regarding the Euclidean distance, but may be dissimilar regarding the shortest path distance (network distance). Thus, for every pair of points between these two trajectory parts, the Euclidean distance is small, but the corresponding network distance is large because the long edges must be crossed. Therefore, it is important in network applications to use the network distance metric instead of the Euclidean metric.

Let $c(v_i, v_j)$ denote the cost function to travel from a source node v_i to a destination node v_j . As we have already mentioned, this cost for the most network-based applications is defined as the shortest path distance (network distance) between the two nodes. In this paper we fix this cost to be the network distance. We also fix the following requirements for the graph representation of the network *G*: *G* must be a directed or non-directed, positive

Fig. 4. Trajectory proximity.

weighted and strongly connected graph. These cases represent successfully the most real network applications (road networks, etc.).

The cost function (network distance) satisfies the following properties:

Property I. The cost function $c(v_i, v_j)$ gives zero values if and only if $v_i \equiv v_j$.

It is obvious that c(v, v) = 0 for any node v in the graph representation. It also holds that $c(v_i, v_j) = 0 \Rightarrow v_i \equiv v_j$, because it has been assumed that any transition between nodes comes at a non-zero cost (positive weighted graphs).

Property II. The cost function $c(v_i, v_j)$, definitely satisfies the positivity property and the triangular inequality:

- $c(v_i, v_j) \ge 0;$
- $c(v_i, v_j) \leq c(v_i, v_x) + c(v_x, v_j).$

Property III. The cost function $c(v_i, v_j)$, does not satisfy in general the symmetric property, therefore it is not definitely a metric function:

• $c(v_i, v_j) \neq c(v_j, v_i)$.

But how does this reflect reality? Consider a directed road network with many one-way road segments, which is quite common. Then, it is clear that if a car goes from a source node v_i to a destination node v_j , it will cover a distance generally different than its way back from v_j to v_i , as it has to pass through different nodes with different weights.

3.1.1. Network distance measure 1

The first network distance measure D_{net1} that we propose uses network-based computations. The distance $d(v_i, v_j)$ between any two nodes v_i and v_j , belonging to trajectories T_a and T_b , respectively, is given by the following definition.

Definition 1. The distance $d(v_i, v_j)$ between two graph nodes v_i and v_j is defined as follows:

$$d(v_i, v_j) = \begin{cases} 0, & \text{if } v_i = v_j, \\ \frac{c(v_i, v_j) + c(v_j, v_i)}{2D_c}, & \text{otherwise,} \end{cases}$$
(2)

where $D_G = \max\{c(v_i, v_j), \forall v_i, v_j \in V(G)\}$ is the diameter of the graph *G* of the spatial network and is a global constant for the applications. Its value can be computed taking the overall maximum of possible values of the cost function.

Proposition 1. The distance function $d(v_i, v_j)$ assumes values in the interval [0, 1].

Proof. This is obvious when the function returns a zero value. Otherwise it returns the ratio $\frac{c(v_i, v_j)+c(v_j, v_i)}{2D_c}$. But, clearly we have: $c(v_i, v_j) \leq D_G$ and $c(v_j, v_i) \leq D_G$, and by summation we get: $c(v_i, v_j) + c(v_j, v_i) \leq 2D_G$. Therefore, by division we get: $d(v_i, v_j) = \frac{c(v_i, v_j)+c(v_j, v_i)}{2D_G} \leq 1$. In addition, we have always $c(v_i, v_j) \geq 0$ and $c(v_j, v_i) \geq 0$ (positivity), thus $d(v_i, v_j) \geq 0$. \Box

Proposition 2. The distance function $d(v_i, v_j)$ satisfies the metric properties.

Proof. We need to prove the following properties for every graph nodes v_i , v_j , v_x :

(i) $d(v_i, v_j) \ge 0;$ (ii) $d(v_i, v_j) = d(v_j, v_i);$ (iii) $d(v_i, v_j) \le d(v_i, v_x) + d(v_x, v_j).$ Clearly, property (i) is true by Proposition 1. Property (ii) is always true if $v_i = v_j$. Otherwise, if $v_i \neq v_j$, we have:

$$d(v_i, v_j) = \frac{c(v_i, v_j) + c(v_j, v_i)}{2D_G} = \frac{c(v_j, v_i) + c(v_i, v_j)}{2D_G} = d(v_j, v_i)$$

Thus, it is true in any case.

Property (iii) is obvious if $v_i = v_j$ or $v_i = v_x$ or $v_j = v_x$. Otherwise, if $v_i \neq v_i \neq v_x$ by substitution we get:

$$\frac{c(v_i, v_j) + c(v_j, v_i)}{2D_G} \leqslant \frac{c(v_i, v_x) + c(v_x, v_i)}{2D_G} + \frac{c(v_x, v_j) + c(v_j, v_x)}{2D_G}$$
(3)

Due to the fact that the cost function satisfies the triangular inequality, we have:

$$c(v_i, v_j) \leq c(v_i, v_x) + c(v_x, v_j)$$

$$c(v_i, v_i) \leq c(v_i, v_x) + c(v_x, v_i)$$

By summation and by division with $2D_G$ we take inequality (3), thus property (iii) has been proven. \Box

Definition 2. The network distance $D_{net1}(T_a, T_b)$ between two trajectories T_a and T_b of description length m is defined as follows:

$$D_{net1}(T_a, T_b) = \frac{1}{m} \sum_{i=1}^{m} (d(v_{ai}, v_{bi}))$$
(4)

Proposition 3. *The distance measure* $D_{net1}(T_a, T_b)$ *assumes values in the interval* [0, 1].

Proof. Omitted.

Proposition 4. The distance measure $D_{net1}(T_a, T_b)$ satisfy the metric properties.

Proof. We need to prove the following properties for every trajectories T_a , T_b , T_x of description length m:

(i) $D_{net1}(T_a, T_b) \ge 0;$ (ii) $D_{net1}(T_a, T_b) = D_{net1}(T_b, T_a);$ (iii) $D_{net1}(T_a, T_b) \le D_{net1}(T_a, T_x) + D_{net1}(T_x, T_b).$

Clearly, property (i) is true by consulting Proposition 3. Property (ii) is true because is also true for the distance function d (Proposition 2), so:

$$D_{net}(T_a, T_b) = \frac{1}{m} \sum_{i=1}^{m} (d(v_{ai}, v_{bi})) = \frac{1}{m} \sum_{i=1}^{m} (d(v_{bi}, v_{ai})) = D_{net}(T_b, T_a)$$

Property (iii) is written equally by substitution:

$$\frac{1}{m}\sum_{i=1}^{m}(d(v_{ai}, v_{bi})) \leqslant \frac{1}{m}\sum_{i=1}^{m}(d(v_{ai}, v_{xi})) + \frac{1}{m}\sum_{i=1}^{m}(d(v_{xi}, v_{bi}))$$
$$\iff \sum_{i=1}^{m}(d(v_{ai}, v_{bi})) \leqslant \sum_{i=1}^{m}(d(v_{ai}, v_{xi})) + \sum_{i=1}^{m}(d(v_{xi}, v_{bi}))$$
(5)

From Proposition 2, we have the following inequalities:

$$d(v_{ai}, v_{bi}) \leq d(v_{ai}, v_{xi}) + d(v_{xi}, v_{bi}) \quad \forall i \in \{1, 2, \dots, m\}$$

By summation we get (5). \Box

Fig. 5 shows two trajectories T_a , T_b for which we are interested to calculate their distance. Assuming that $D_G = 100$, we have the following calculations:

$$d(v_{ai}, v_{bi}) = \left\{\frac{17}{200}, \frac{16}{200}, \frac{9}{200}, \frac{7}{200}, 0, \frac{5}{200}, \frac{13}{200}\right\}$$

E. Tiakas et al./The Journal of Systems and Software xxx (2009) xxx-xxx



Fig. 5. Trajectory similarity example.

$$D_{net1}(T_a, T_b) = \frac{1}{7} \left(\frac{17}{200} + \frac{16}{200} + \frac{9}{200} + \frac{7}{200} + 0 + \frac{5}{200} + \frac{13}{200} \right)$$
$$= \frac{1}{7} \frac{67}{200} = 0.047857$$

3.1.2. Network distance measure 2

The second distance measure, D_{net2} , that we propose uses an Euclidean-based distance function (*de*) in combination with the previous global constant D_G (the graph diameter by the network distance).

It can be used for fast calculations only for graphs where the coordinates of the nodes are available. In fact, in many cases the Euclidean distance results in poor performance regarding the quality of results. However, as it will be described later, it offers a "quick-and-dirty" view of the results.

Definition 3. The distance $de(v_i, v_j)$ between two graph nodes v_i and v_i is defined as follows:

$$de(v_i, v_j) = \frac{euclidean(v_i, v_j)}{D_G} = \frac{\sqrt{(x_{v_i} - x_{v_j})^2 + (y_{v_i} - y_{v_j})^2}}{D_G}$$
(6)

where x_{v_i} , y_{v_i} are the coordinates of node v_i , and x_{v_j} , y_{v_j} are the coordinates of node v_i .

Proposition 5. The distance function $de(v_i, v_j)$ assumes values in the interval [0, 1].

Proof. Let DE_G be the maximum Euclidean distance between all nodes of the graph representing the spatial network: $DE_G = \max\{euclidean(v_i, v_j), \forall v_i, v_j \in V(G)\}$. Then it is obvious that:

$$euclidean(v_i, v_i) \leq DE_G \leq D_G \quad \forall v_i, v_i \in V(G)$$

The last inequality holds because all network distances are always greater than or equal to the corresponding Euclidean distances. Therefore, we have:

$$\frac{euclidean(v_i, v_j)}{D_G} \leqslant 1 \iff de(v_i, v_j) \leqslant 1$$

Moreover, as all distances are positive (or zero when $v_i = v_j$), we have always: $de(v_i, v_j) \ge 0$. \Box

Proposition 6. The distance function $de(v_i, v_j)$ satisfies the metric properties.

Proof. Due to the fact that the Euclidean distance *euclidean*(v_i, v_j) satisfies the metric properties and $de(v_i, v_j)$ is the Euclidean distance divided by the positive constant D_G , it is evident that $de(v_i, v_j)$ also satisfies the metric properties. \Box

Definition 4. The network distance $D_{net2}(T_a, T_b)$ between two trajectories T_a and T_b of description length m is defined as follows:

$$D_{net2}(T_a, T_b) = \frac{1}{m} \sum_{i=1}^{m} (de(v_{ai}, v_{bi}))$$
(7)

Proposition 7. The distance measure $D_{net2}(T_a, T_b)$ assumes values in the interval [0, 1].

Proof. Omitted.

Proposition 8. The distance measure $D_{net2}(T_a, T_b)$ satisfy the metric properties.

Proof. Omitted.

3.2. Incorporating time information

The similarity measures defined in the previous section take into consideration only the traveling cost information, which depends on the spatial network. In applications such as traffic analysis, the time information associated with each trajectory is also very important.

Definition 5. Given two trajectories $T_a \in \mathcal{T}$ and $T_b \in \mathcal{T}$ of description length *m*, their distance with respect to time $D_{time}(T_a, T_b)$ is given by

$$D_{time}(T_a, T_b) = \frac{1}{m-1} \sum_{i=1}^{m-1} \frac{|(T_a[i+1].t - T_a[i].t) - (T_b[i+1].t - T_b[i].t)|}{\max\{(T_a[i+1].t - T_a[i].t), (T_b[i+1].t - T_b[i].t)\}}$$

Essentially, the time similarity between two trajectories, as it has been defined, measures their resemblance with respect to the time required to travel from one node to the next (inter-arrival times).

Fig. 6 depicts some examples for the time similarity calculations, where we have three trajectory parts T_a , T_b , T_c with the same description length and the inter-arrival times appear next to their directed edges.

With the previous definition we have the following calculations:

$$D_{time}(T_a, T_b) = \frac{1}{4} \left(\frac{1}{5} + \frac{0}{7} + \frac{1}{4} + \frac{0}{2} \right) = 0.1125$$
$$D_{time}(T_a, T_c) = \frac{1}{4} \left(\frac{0}{5} + \frac{4}{7} + \frac{2}{6} + \frac{2}{4} \right) = 0.35119$$



Fig. 6. Time similarity calculation example.

We observe that T_a is more similar to T_b than T_c and this happens because the corresponding inter-arrival times of the pair T_a , T_b are much closer.

Proposition 9. The distance measure $D_{time}(T_a, T_b)$ assumes values in the interval [0, 1].

Proof. Omitted.

Proposition 10. *The distance measure* $D_{time}(T_a, T_b)$ *satisfy the metric properties.*

Proof. We need to prove the following properties for any trajectories T_a , T_b , T_x of description length m:

- (i) $D_{time}(T_a, T_b) \ge 0;$
- (ii) $D_{time}(T_a, T_b) = D_{time}(T_b, T_a);$
- (iii) $D_{time}(T_a, T_b) \leq D_{time}(T_a, T_x) + D_{time}(T_x, T_b)$.

Clearly, property (i) is true by Proposition 9. Let us denote the inter-arrival times of all trajectory parts of T_a , T_b and T_x as follows: $\delta_{ai} = T_a[i+1].t - T_a[i].t$, $\delta_{bi} = T_b[i+1].t - T_b[i].t$ and $\delta_{xi} = T_x[i+1].t - T_x[i].t$, for all i = 1, 2, ..., m-1. Then, property (ii) is true because we have:

$$D_{time}(T_a, T_b) = \frac{1}{m-1} \sum_{i=1}^{m-1} \frac{|\delta_{ai} - \delta_{bi}|}{\max\{\delta_{ai}, \delta_{bi}\}} = \frac{1}{m-1} \sum_{i=1}^{m-1} \frac{|\delta_{bi} - \delta_{ai}|}{\max\{\delta_{bi}, \delta_{ai}\}} = D_{time}(T_b, T_a)$$

By substitution, property (iii) is written as

$$\frac{1}{m-1} \sum_{i=1}^{m-1} \frac{|\delta_{ai} - \delta_{bi}|}{\max\{\delta_{ai}, \delta_{bi}\}} \leqslant \frac{1}{m-1} \sum_{i=1}^{m-1} \frac{|\delta_{ai} - \delta_{xi}|}{\max\{\delta_{ai}, \delta_{xi}\}} \\
+ \frac{1}{m-1} \sum_{i=1}^{m-1} \frac{|\delta_{xi} - \delta_{bi}|}{\max\{\delta_{xi}, \delta_{bi}\}} \iff \sum_{i=1}^{m-1} \frac{|\delta_{ai} - \delta_{bi}|}{\max\{\delta_{ai}, \delta_{bi}\}} \\
\leqslant \sum_{i=1}^{m-1} \frac{|\delta_{ai} - \delta_{xi}|}{\max\{\delta_{ai}, \delta_{xi}\}} + \sum_{i=1}^{m-1} \frac{|\delta_{xi} - \delta_{bi}|}{\max\{\delta_{xi}, \delta_{bi}\}}$$
(8)

It is sufficient to prove the following inequalities $\forall i = 1, ..., m - 1$:

$$\frac{|\delta_{ai} - \delta_{bi}|}{\max\{\delta_{ai}, \delta_{bi}\}} \leqslant \frac{|\delta_{ai} - \delta_{xi}|}{\max\{\delta_{ai}, \delta_{xi}\}} + \frac{|\delta_{xi} - \delta_{bi}|}{\max\{\delta_{xi}, \delta_{bi}\}}$$
(9)

To prove (9) it is enough to prove that for every positive numbers *a*, *b*, *c* the following inequality holds:

$$\frac{|a-b|}{\max\{a,b\}} \leqslant \frac{|a-c|}{\max\{a,c\}} + \frac{|c-b|}{\max\{c,b\}}$$
(10)

But, this inequality is obvious if a = b, or a = c, or b = c, and for all other ordering cases of the numbers a, b, c also holds:

$$\frac{b-a}{b} \leq \frac{c-a}{c} + \frac{c-b}{c}$$
$$\iff c(b-a) \leq b(c-a) + b(c-b)$$
$$\iff (b+a)(c-b) \geq 0$$

which it holds as a, b are positive and b < c. • If a < c < b then it gives:

$$\frac{b-a}{b} \leq \frac{c-a}{c} + \frac{b-c}{b}$$
$$\iff c(b-a) \leq b(c-a) + c(b-c)$$
$$\iff (c-a)(b-c) \geq 0$$

which it holds as a < c and c < b. • If b < a < c then it gives:

$$\frac{a-b}{a} \leqslant \frac{c-a}{c} + \frac{c-b}{c}$$
$$\iff c(a-b) \leqslant a(c-a) + a(c-b)$$
$$\iff (a+b)(c-a) \ge 0$$

which it holds as a, b are positive and a < c. • If b < c < a then it gives:

$$\frac{a-b}{a} \leq \frac{a-c}{a} + \frac{c-b}{c}$$
$$\iff c(a-b) \leq c(a-c) + a(c-b)$$
$$\iff (c-b)(a-c) \ge 0$$

which it holds as b < c and c < a. • If c < a < b then it gives:

$$\frac{b-a}{b} \leq \frac{a-c}{a} + \frac{b-c}{b}$$
$$\iff a(b-a) \leq b(a-c) + a(b-c)$$
$$\iff (a+b)(a-c) \geq 0$$

which it holds as a, b are positive and c < a.

• If
$$c < b < a$$
 then it gives:

$$\frac{a-b}{a} \leq \frac{a-c}{a} + \frac{b-c}{b}$$

$$\iff b(a-b) \leq b(a-c) + a(b-c)$$

$$\iff (a+b)(b-c) \ge 0$$

which it holds as a, b are positive and c < b.

Therefore, inequality (10) is true, and property (iii) has been proven. $\ \Box$

3.2.1. Spatio-temporal similarity measures and methods

We have at hand different distance measures, D_{net} and D_{time} , that can be used to compare trajectories of the same length in space and time. Several applications may require both similarity measures to extract useful knowledge.

There are three different methods in order to retrieve similar trajectories in space-time as proposed in Hwang et al. (2005): (i) searching similar trajectories with direct application of spatio-temporal distance measures, (ii) filtering trajectories based on temporal similarity and refining similar trajectories based on spatial distance, (iii) filtering trajectories based on spatial distance, (iii) filtering trajectories based on spatial distance, (iii) filtering trajectories based on temporal similar trajectories based on temporal similarity and refining similar trajectories based on temporal distance.

Here we suggest the methods (i) and (iii), due to the fact that method (ii) can hardly be found in practical applications.

To implement method (i) we can combine the two distance measures D_{net} and D_{time} into a single one. For example, the two distances may be weighted with parameters W_{net} and W_{time} such that $W_{net}+W_{time}=1$. The total (combined) distance can then be expressed as follows:

$$D_{total}(T_a, T_b) = W_{net} \cdot D_{net}(T_a, T_b) + W_{time} \cdot D_{time}(T_a, T_b)$$

It is evident that the distance measure D_{total} satisfies the metric space properties. However, this approach poses a significant limitation, since the values of W_{net} and W_{time} must be known in advance.

Consequently we propose method (iii) using D_{net} and D_{time} separately, where the distance D_{net} making the filtering step in space and the distance D_{time} making the refinement step in time. In this way, two parameter distances are required to be posed by the query. The distance E_{net} expresses the desired similarity with respect to the D_{net} distance measure, whereas the distance E_{time} expresses the desired similarity regarding the D_{time} distance measure.

If the user wishes to focus only on the network distance, then the value of E_{time} may be set to 1. Otherwise, another value is required for E_{time} , which determines the desired similarity in the time domain. By allowing the user to control the values of E_{net} and E_{time} a significant degree of flexibility is achieved, since the "weight" of each distance can be controlled at will.

4. Indexing and query processing issues

In this section, we study some important issues regarding trajectory similarity. Firstly, we discuss the problem of handling trajectories of different description length, by decomposing a trajectory to sub-trajectories. Then, we study the use of indexing schemes for sub-trajectories. Finally, we study some fundamental query processing issues.

4.1. Trajectory decomposition

Up to now we have handled the case where all trajectories are of the same description length. We proceed now to relax this assumption, by considering trajectories of different lengths. In fact, this is the more general case that reflects reality. First of all, two trajectories may involve a different number of visited nodes, and therefore their description length will be different. Furthermore, we cannot always guarantee that moving objects report their positions at fixed time intervals. Due to noise, several measurements may be lost, or different moving objects report their positions at different time intervals. In these cases, two trajectories may have different description lengths.

Let *T* be a trajectory of description length *m*. Moreover, let μ denote an integer such that $\mu \leq m$. *T* is decomposed into $m - \mu + 1$ sub-trajectories, by using a window of length μ , and progressively moving one node at a time from left to right. Each of the resulting sub-trajectories has a length of μ . Fig. 7 illustrates an example of the decomposition process, where m = 6 and $\mu = 3$.

By following the same process for all trajectories $T \in \mathcal{T}$ we get a new set of sub-trajectories S, all of description length μ . Moreover, we have already defined a distance measure for trajectories of the



same description length in the previous section given by either D_{net} or D_{time} which both satisfy the metric space properties.

4.2. Indexing schemes

Our next step involves indexing the set S of sub-trajectories, enabling efficient query processing. Towards this direction, we propose two schemes, which are both based on the M-tree access method (Ciaccia et al., 1997). Note that since a vector representation of each sub-trajectory is not available, techniques like R-trees (Guttman, 1984) and its variants are not applicable. Recall that, the M-tree is already equipped by the necessary tools to handle range and nearest-neighbor queries, as it has been reported in Ciaccia et al. (1997). The only requirement for the M-tree to work properly is that the distance used must satisfy the metric space properties. Since both D_{net} and D_{time} satisfy these properties, they can be used as distance measures in M-trees. Note that, among the metric indexing schemes we choose the M-tree because of its simplicity. However, other secondary memory schemes for metric spaces or any other metric access method can been applied equally well (e.g., SlimTrees; Traina et al., 2000). Two alternatives are followed towards indexing sub-trajectories:

- *M-treel method*. In this scheme, only the NET-M-tree is used to check the constraint regarding *E*_{net}. Then, in a subsequent step the candidate sub-trajectories are checked against the time constraints. This way, only one M-tree is used.
- *M*-treell method. In this scheme, two M-trees are used to handle D_{net} and D_{time} separately. These trees are termed NET-M-tree and TIME-M-tree, respectively. Each M-tree is searched separately using E_{net} and E_{time} , respectively. Then, the intersection of both results is determined to get the sub-trajectories that satisfy the network and time constraints.

4.3. Query processing fundamentals

A user query is defined by a triplet $\langle T_q, E_{net}, E_{time} \rangle$ where T_q is the query trajectory, E_{net} is the radius for the network distance and E_{time} is the radius for the time distance. For the query processing to be consistent with the proposed framework, each query trajectory T_q must be of at least description length μ . If this is not true, padding is performed by repeating, for example, the last node of the trajectory several times, until the description length μ is reached. In the general case where the description length of T_q is greater than μ , the decomposition process is applied to obtain the sub-trajectories of T_q . Finally, if the description length of T_q is equal to μ , then only one sub-trajectory is produced.

Let p denote the number of sub-trajectories of T_q determined by the trajectory decomposition process. The next step depends on the indexing scheme we utilize, i.e. either M-treel or M-treell as they have been described previously. A trajectory is part of the answer if there is at least one of its sub-trajectories that satisfy the network and time constraints for at least one query sub-trajectory. In the sequel, we analyze the whole process in detail:

- Having a query trajectory T_q of description length l and the E_{net} , E_{time} parameters, we decompose T_q into $p = l \mu + 1$ sub-trajectories (if $l > \mu$) with the window method and then we construct their set $QS(T_q)$.
- For every query sub-trajectory $qs \in QS(T_q)$, we execute a simple range query to NET-M-Tree with radius E_{net} and collect related sub-trajectories into the set C_{net} .
- If M-treell method is used then we execute another simple range query to TIME-M-Tree with radius *E*_{time} and collect related sub-trajectories into the set *C*_{time}.

1- -

8

- If M-treel method is used then we check every sub-trajectory in C_{net} against E_{time} and from the selected results we construct the set AS. Otherwise, If M-treelI method is used, the results' set AS is constructed with the common sub-trajectories of the sets C_{net} and C_{time}. In both cases, the set AS contains the resulted sub-trajectories ID's.
- From the set *AS* we take the corresponding trajectories ID's and we construct the final result set *AT*.

In any case, a trajectory $T \in \mathcal{T}$ will appear in the result set, if and only if there exists at least one sub-trajectory *ts* of *T* which is similar to at least one sub-trajectory *qs* of the query trajectory T_q , and also satisfies the network and time constraints. More formally:

T is similar to $T_q \iff \exists ts \subseteq T, \exists qs \subseteq T_q : D_{net}(ts, qs)$

 $\leq E_{net} \wedge D_{time}(ts, qs) \leq E_{time}$

Fig. 8 presents an outline of the algorithm. Taking into account that the consecutive sub-trajectories of T_q have $(\mu - 1)$ common nodes, most calculations and requests can be already in the memory, as we check one sub-trajectory after another, so it is strongly recommended to use an LRU memory buffer.

4.4. Distance buffering

The distance measure D_{net1} uses the shortest path distance between graph nodes. These computations can be performed more efficiently by using an LRU buffer. The LRU buffer maintains a constant amount of distance values into main memory. In the experimental results section we show that only a relatively small buffer size is adequate to accelerate performance, offering a good hit ratio. If the network graph has at most a few thousand nodes, it is suggested to precompute all distances $c(v_i, v_j)$ between nodes and to put them into a hash-based file. Then, the LRU memory buffer can cooperate with this file during the request procedure for even better performance. Later, we discuss the alternative of storing only a subset of precomputed distances on the disk, to handle large graphs.

The algorithm in Fig. 9 illustrates the process of retrieving a distance $c(v_i, v_j)$. The variables *requests*, *hits*, and *misses* are used to test buffer performance.

It is important to remind that the LRU memory buffer and the precomputed distances disk file, are used only with D_{net1} . They are not necessary for D_{time} calculations and in D_{net2} measure which does not use network distances at all.

4.5. Combining measures D_{net1} and D_{net2} (filtering and refinement)

Due to network restrictions, a similarity range query using the D_{net2} distance measure may return some trajectories that are not similar regarding distance measure D_{net1} (false alarms). This effect is more significant when the shortest path distance between nodes is considerably higher than their Euclidean distance. Therefore, we need to detect these trajectories using another measure, which respects the network restrictions in space, and use it in a refinement step during query processing. For this reason, we can select the distance measure D_{net1} to handle false alarm detection. This procedure will give correct results if and only if we prove that every trajectory that appears in the result set of D_{net1} measure, appears also in the result set of D_{net2} , when we apply an E_{net} range query.

Proposition 11. For every two trajectories T_a , T_b the following inequality always holds:

Algorithm SimilaritySearch $(T_q, E_{net}, E_{time}, \mu)$				
Input				
T_q : query trajectory				
E_{net} : network distance radius				
E_{time} : time distance radius				
μ : minimum description length of query sub-trajectory				
Output				
AS: set of sub-trajectory IDs				
AT: set of trajectory IDs				
1. $QS(T_q) = $ all sub-trajectories of T_q of description length μ				
2. for each query sub-trajectory $qs \in QS(T_q)$				
3. if method M-treeI is used then				
4. search NET-M-tree using qs and E_{net}				
5. $C_{net} = \text{candidate sub-trajectories from NET-M-tree}$				
6. check every sub-trajectory in C_{net} against E_{time}				
7. update AS				
8. else if method M-treeII is used then				
9. search NET-M-tree using qs and E_{net}				
10. $C_{net} = \text{candidate sub-trajectories from NET-M-tree}$				
11. search TIME-M-tree using qs and E_{time}				
12. $C_{time} = \text{candidate sub-trajectories from TIME-M-tree}$				
13. $AS = C_{net} \cap C_{time}$				
14. end if				
15. end for				
16. calculate AT from AS				
17. $\mathbf{return}(AS,AT)$				

Fig. 8. Outline of similarity search algorithm.

E. Tiakas et al./The Journal of Systems and Software xxx (2009) xxx-xxx

Algorithm RetrieveDistance (v_i, v_j)
Input
v_i : source node
v_j : destination node
Output
$c(v_i, v_j)$: value of the cost function between nodes v_i and v_j

1. $requests++$				
2. search in LRU memory buffer for distance $c(v_i, v_j)$				
3. if distance found in buffer then				
4. $\mathbf{return}(c(v_i, v_j))$				
5. $hits++$				
6. else				
7. if a precomputed distances disk file is used then				
8. open disk file				
9. find record with distance $c(v_i, v_j)$				
10. $\mathbf{return}(c(v_i, v_j))$				
11. insert distance $c(v_i, v_j)$ in memory buffer with LRU rule				
12. $misses++$				
13. else				
14. compute the distance $c(v_i, v_j)$				
15. $\mathbf{return}(c(v_i, v_j))$				
16. insert distance $c(v_i, v_j)$ in memory buffer with LRU rule				
17. $misses++$				
18. end if				
19. end if				

Fig. 9. Outline of distance retrieval algorithm.

 $D_{net2}(T_a, T_b) \leq D_{net1}(T_a, T_b)$

Proof. As the shortest path distance $c(v_i, v_j)$ between two graph nodes v_i , v_j is always greater than or equal to their corresponding Euclidean distance, it always holds that:

 $euclidean(v_i, v_j) \leq c(v_i, v_j) \quad \forall v_i, v_j \in V$

By dividing with the constant D_G we get:

$$\frac{euclidean(v_i, v_j)}{D_G} \leqslant \frac{c(v_i, v_j)}{D_G} \iff de(v_i, v_j) \leqslant d(v_i, v_j) \quad \forall v_i, \ v_j \in V$$

Therefore, for every two trajectories $T_a = (v_{a1}, v_{a2}, ..., v_{am})$ and $T_b = (v_{b1}, v_{b2}, ..., v_{bm})$, where $v_{ai} \in V$ and $v_{bi} \in V$ ($\forall i = 1, ..., m$), we have the following inequalities:

$$de(v_{ai}, v_{bi}) \leqslant d(v_{ai}, v_{bi}) \quad \forall i = 1, \dots, m$$

By summation, we get:

$$\sum_{i=1}^{m} (de(v_{ai}, v_{bi})) \leq \sum_{i=1}^{m} (d(v_{ai}, v_{bi})) \iff \frac{1}{m} \sum_{i=1}^{m} (de(v_{ai}, v_{bi}))$$
$$\leq \frac{1}{m} \sum_{i=1}^{m} (d(v_{ai}, v_{bi})) \iff D_{net2}(T_a, T_b) \leq D_{net1}(T_a, T_b)$$

and the proposition has been proven. \Box

Following Proposition 11, when we have a query trajectory T_q and a network query range E_{net} , all trajectories returned by D_{net1} measure will appear in the result set of D_{net2} , because:

$$D_{net2}(T_q, T) \leqslant D_{net1}(T_q, T) \leqslant E_{net} \quad \forall T \in \mathcal{T}$$

Fig. 10 illustrates the outline of the similarity search algorithm including the refinement step. D_{net2} is used as the filtering distance measure, whereas D_{net1} is used for refinement, to eliminate false

alarms. An important observation is that this scheme can be applied to both M-treel and M-treell methods, and moreover, it can be used with any well-defined distance measure, as long as the following lower-bounding property holds:

$$D_{\text{filtering}}(T_a, T_b) \leqslant D_{\text{refinement}}(T_a, T_b) \quad \forall T_a, \ T_b \in \mathcal{T}$$

5. Performance evaluation

In this section, we give information about the implementation of the proposed approach in C++ and the results of experiments that confirm and evaluate all previous algorithms, procedures and techniques. All experiments have been conducted on a Pentium IV running Windows XP, with 1 GB of RAM, and a 320 GB-SATA2-16 MB hard disk. First, we present the construction of used spatial network and trajectory data set. Then, we present the construction of M-Trees for each defined measure and how the proposed measures express well the notion of similarity in space and time. At the main part, we present the evaluation results of all proposed methods for similarity range queries.

5.1. Spatial network data

All experiments have been conducted using a real-world spatial network, the road network of Oldenburg city (Brinkhoff, 2002). The cost function $c(v_i, v_j)$ between two nodes of the graph representation is the shortest path distance. The number of vertices in the Oldenburg data set is 6105. Therefore, the total number of precomputed distances among all possible pairs of vertices is 37,271,025. These distances are stored in a hash-based file on disk (DISTfile), using the Hilbert space filling curve as a hashing function. The Hilbert curve values are derived from the corresponding source and target node ID's of the distances, which are integers into the

E. Tiakas et al./The Journal of Systems and Software xxx (2009) xxx-xxx

Algorithm SimilaritySearchWithRefinement($T_q, E_{net}, E_{time}, \mu$) Input T_q : query trajectory E_{net} : network distance radius E_{time} : time distance radius μ : minimum description length of query sub-trajectory Output ASF: final set of sub-trajectory IDs ATF: final set of trajectory IDs 1. $QS(T_q) =$ all sub-trajectories of T_q of description length μ 2. $ASF = \emptyset$ 3. for each query sub-trajectory $qs \in QS(T_q)$ if method M-treeI is used 4. 5.search NET-M-tree (constructed by the basic metric) using qs and E_{net} 6. C_{net} = candidate sub-trajectories from NET-M-tree (using the D_{net} distance of the basic metric) 7. check every sub-trajectory in C_{net} against E_{time} update AS8. else if method M-treeII is used then 9. 10. search NET-M-tree (constructed by the basic metric) using qs and E_{net} 11. C_{net} = candidate sub-trajectories from NET-M-tree (using the D_{net} distance of the basic metric) 12.search TIME-M-tree using qs and E_{time} $C_{time} =$ candidate sub-trajectories from TIME-M-tree 13.14. $AS = C_{net} \cap C_{time}$ 15.end if 16.for each sub-trajectory S_i in AS17.compute the D_{net} distance of S_i from qs using the selected refinement metric 18. insert S_i in ASF if that distance is less than or equal to E_{net} 19. end for 20. end for 21. calculate ATF from ASF22. return(ASF, ATF)

Fig. 10. Outline of similarity search algorithm with refinement step.

interval $[0, |V_G| - 1]$, (e.g., for the distance $c(v_i, v_j)$ the value $Hilbert(ID(v_i), ID(v_j))$ is calculated). For the selected road network, the total time required for all precomputations and creation of DISTfile is 3,180.581 s. The record length has been set to 16 bytes, so the final file capacity is 596,336,400 bytes (285 MB zipped).

An in-core LRU buffer has been used to keep a number of precomputed distances in main memory (we initialized the buffer selecting some top-used distances through calculations which actually are distances between nodes that included in the most trajectory parts). The size of the buffer has been set to 2000, which is a relatively small value compared to the total number of pair-wise distances. We have computed the average number of network distance calculation requests, the average number of hits and misses, in simple range queries in space using D_{net1} and D_{net2} . The results show that almost 85% of the distance requests are absorbed by the main memory buffer and therefore, we avoid fetching them from the disk. The more buffer pages are available, the higher the hit ratio becomes.

The fast retrieval of shortest path distances is the most time consuming factor affecting the performance of network-based distance calculations, the construction of M-Trees and finally in the performance of similarity range queries.

5.2. Construction of trajectories and sub-trajectories

The trajectory data set T we have used for the experiments consists of 3797 trajectories of objects moving on the road segments of

Oldenburg city, using the generator developed in Brinkhoff (2002). Each trajectory has a minimum description length of 10 and maximum description length of 100 nodes. A sliding window of description length $\mu = 10$ has been used to generate the sub-trajectories of each trajectory. Therefore, the total number of sub-trajectories produced (set S) is 75,144.

Moreover, it is important to study the distribution of the constructed trajectory data set among the nodes of the road network. This will help to evaluate if the data set represents well a realworld trajectory set of this town. So, we record in a new file all node ID's used by the trajectories, with the frequency that are being used (how many trajectories pass through) in a descending order. Fig. 11a gives the recorded distribution and Fig. 11b depicts the top-100 most used nodes in the network by the trajectories.

It is evident that we have a skew distribution of nodes in trajectories and this reflects reality: there are some nodes that are being used very often which are center points of this town or hard traffic points, and the most peripheral nodes are being used much rarely. Therefore, our trajectory data set is a good representative of a real traffic condition.

5.3. M-tree construction

We have constructed four different M-trees. The NET-M-trees which are implemented based on the $D_{net1,2,3}$ measures and the TIME-M-tree implemented based on the D_{time} measure. Recall that,

E. Tiakas et al. / The Journal of Systems and Software xxx (2009) xxx-xxx



Fig. 11. (a) Distribution of nodes in trajectories; (b) top-100 most frequent nodes.

all M-trees handle the same set S of sub-trajectories of description length $\mu = 10$ and not the complete trajectories of moving objects.

We have utilized the bulk-loading method for the construction of all M-trees, and the following parameters values have been used: a page size of 4 KB, 5% minimum node utilization, minimum overlaps promote part and root functions, a general hyper-plane split strategy, and radius function by average. Table 2 shows the total number of network distances computed during the construction, the number of zero distances, the final file capacity of Mtrees on disk and the total construction time. Note that we have exploited precomputed distances (LRU buffer and DISTfile) during the construction procedure.

We observe that D_{net2} gives the smallest capacity and construction time, because network distance computations are not required.

5.4. Evaluation of similarity measures

We have randomly selected several trajectories from different areas of Oldenburg and we have performed similarity range queries by using all measures.

Figs. 12 and 13 show the results of range queries with radius $E_{net} = 0.01, 0.05, 0.10$, in a random selected query trajectory from our data set, using the available network distance measures. By studying these figures we observe that:

- In all metrics, the resulted trajectories firstly appeared in the closest neighbor of query trajectory and as the radius *E_{net}* increases, they expand into connected and almost rounded areas, in which the query trajectory takes a central position.
- All query trajectory results of D_{net1} metric (using a constant E_{net} range) are included in the results of D_{net2} metric, according to Proposition 11.

Table 2	
Information regarding the construction of M-trees.	

M-tree	Distances	Zeros	Capacity (MB)	Time
D _{net1}	1,574,890	38,309	32.5	13 min + 7 s
D _{net2}	1,494,416	37,761	32.1	35 s
D _{time}	4,013,864	40,461	30.9	1 min + 46 s

5.5. Performance evaluation of M-treeI and M-treeII methods

In this section, we study the performance of M-treel,II methods using all the proposed metrics. We selected randomly 100 trajectories from our data set and from different parts of the town and we performed similarity range queries using the M-treel,II methods. We gave all combination values into the interval [0, 1] with a step of 0.05 in E_{net} , E_{time} parameters. The final reported results correspond to the average values of these 100 queries. The basic parameters that are studied are summarized in Table 3.

Fig. 14a depicts the number of similar sub-trajectories found using all available network-based distance measures. Recall, that the results are the same for both M-treel and M-treell methods. As the E_{net} radius increases, D_{net2} first reaches the upper limit (75,144), followed by D_{net1} . Evidently, the distance measure D_{net2} gives more results than D_{net1} due to the lower-bounding property.

Fig. 14b depicts the total time spent for network-based computations using all network-based distance measures. It is evident that D_{net2} is the less time consuming measure since distances are computed by using the Euclidean distance of the nodes. The results are similar for both M-treel and M-treell methods.

Fig. 15 illustrates the memory LRU buffer activity using D_{net1} . Note that D_{net2} does not use the LRU buffer, since no network distance computations are performed. In both cases, the total number of distance hits is about 85%, so with only 2000 buffer pages we have a satisfactory hit ratio. Again, the results are similar for M-treel and M-treell methods.

Fig. 16 depicts the number of time-based distance computations for M-treel and M-treell methods. We observe that in M-tree-II method the number of time-based distance computations depends only on the E_{time} radius, whereas in M-treel method this value depends on both E_{net} and E_{time} parameters, because in this case the total number of computed distances is equal to the number of sub-trajectory results returned by the NET-M-tree. Therefore, we expect less time-based computations in the M-treel method than in the M-treell method. This results in slightly better performance for the M-treel method regarding time-based distance computation overhead, as it is illustrated in Fig. 17.

Fig. 18 depicts the percentage of false alarms for D_{net2+1} , for various values of parameters E_{net} and E_{time} . In the left part of the figure, these parameters change freely, whereas in the right part always

Please cite this article in press as: Tiakas, E. et al., Searching for similar trajectories in spatial networks, J. Syst. Software (2009), doi:10.1016/j.jss.2008.11.832

12

E. Tiakas et al. / The Journal of Systems and Software xxx (2009) xxx-xxx



Fig. 12. Preview of range queries using distance measure *D*_{net1}.



Fig. 13. Preview of range queries using distance measure *D*_{*net*2}.

Table 3

Basic variables measured throughout experiments.

Variable	Description
N _{net}	Number of similar sub-trajectories found in NET-M-tree
D _{net}	Number of network-based distance computations
T _{net}	Total searching time in NET-M-tree (s)
MBFR	Memory LRU buffer total requests
MBFH	Memory LRU buffer total hits
MBFM	Memory LRU buffer total misses
DBFR	Disk LRU buffer total requests
DBFH	Disk LRU buffer total hits
DBFM	Disk LRU buffer total misses
N _{time}	Number of similar sub-trajectories found in TIME-M-tree
D _{time}	Number of time-based distance computations
T _{time}	Total searching time in TIME-M-tree (M-treeII) or in time calculations (M-treeI) (s)
TT	Total query time
AS	Total number of common (M-treeII) or accepted (M-treeI) sub-trajectories found (Net&Time)
AT	Total number of similar trajectories found (final results)
FA	False alarms for sub-trajectories in D_{net2+1} method

 E_{net} equals E_{time} . It is evident, that the existence of false alarms cannot be avoided, due to the distance lower-bounding. However, the percentage of false alarms is relatively small, and therefore effective filtering is performed by applying the Euclidean distance prior to network distance computations. The maximum number of false alarms (around 25%) appears when $E_{net} = 0.25$ and $E_{time} = 0.30$.

Fig. 19a and b depict some representative results regarding the performance of M-treeI and M-treeII methods for $E_{net} = E_{time}$, using all network distance measures. D_{net2} is the most efficient tool but needs validation of correctness, and D_{net2+1} method is the most attractive alternative that can be used for trajectory similarity search, if efficiency is important. However, care should be taken since the usage of D_{net2+1} involves determination of false alarms.

If the number of false alarms is large, performance degradation may appear.

In all the experiments conducted, the method that uses only one M-tree performs marginally better than the method that utilizes two M-trees (one for D_{net} and one for D_{time}). However, the existence of two M-trees offers a higher degree of flexibility during query processing, since we can search for similar trajectories based: (i) only on network distance D_{net} , (ii) only on time distance D_{time} and (iii) both on network and time distances D_{net} and D_{time} . Moreover, different clustering schemes can be applied. More specifically, using the two separate M-trees, a clustering algorithm can provide clusters for D_{net} or D_{time} . Finally, more choices for query optimization are available if both indexes are utilized, since the

E. Tiakas et al. / The Journal of Systems and Software xxx (2009) xxx-xxx



Fig. 14. Number of sub-trajectories (a) and search time (b) for NET-M-tree.



Fig. 15. Memory buffer activity for *D*_{net1}.

query execution engine can form an efficient query execution plan according to the selectivities of the search distances E_{net} and E_{time} , and traverse the M-trees accordingly.

5.6. Impact of precomputed distances

The previous experiments have been conducted by having all network distances precomputed and stored on disk. It has been observed that the precomputation reduces the required computational costs during network-based distance calculations. However, the precomputation assumption may not be realistic in very large spatial networks containing many thousands of nodes. However, even for small spatial networks, if the main memory buffer fails to achieve an acceptable hit ratio, many distance computations will be invoked, resulting in performance degradation.

Figs. 20 and 21 show some interesting results regarding the performance of trajectory similarity queries, when only a subset of the total distances are precomputed. The performance of D_{net1} measure is illustrated in Fig. 20, which depicts the activity of the memorybased (a) as well as the disk-based buffer (b). It is evident that by increasing the number of precomputed distances the total running time of trajectory similarity queries decreases but the cost is still significant, raising problems for ad-hoc query processing. On the other hand, the use of the Euclidean distance for filtering purposes results in a much more efficient scheme, as it is illustrated in Fig. 21.



Fig. 16. Number of time-based distance computations in M-treeI and M-treeII methods.

14

E. Tiakas et al./The Journal of Systems and Software xxx (2009) xxx-xxx



Fig. 17. CPU time (in s) required for time-based distance computations in M-treeI and M-treeII methods.



Fig. 18. Percentage of false alarms for D_{net2+1} method.





E. Tiakas et al. / The Journal of Systems and Software xxx (2009) xxx-xxx



Fig. 20. Memory and disk buffer activity for variable disk buffer sizes.



Fig. 21. Total running time for D_{net2+1} for variable query radius and disk buffer sizes.

6. Conclusions

Although there is significant research work performed on trajectory similarity on moving objects trajectories, the vast majority of the proposed approaches assume that objects can move freely without any motion restrictions. In this paper, we have studied the problem of trajectory similarity query processing in networkconstrained moving objects. We have defined two concepts of similarity. The first is based on the network distance and the second is based on the time characteristics of the trajectories. By using these concepts, we have defined distance measures D_{net} to capture the network similarity and a distance measure D_{time} to capture the time-based similarity of trajectories. All proposed measures satisfy the metric space properties, and therefore, metric-based access methods can be used for efficient indexing and searching.

To support trajectories of different description lengths, a decomposition process is applied. Each trajectory is split to a number of sub-trajectories, which are then indexed by M-trees. The NET-M-tree is used for the D_{net} measure, whereas the TIME-M-tree is used for the D_{time} measure. Two methods have been studied: (i) the M-treel method, which uses only the NET-M-tree and (ii) the M-treel method, which utilizes both trees. Performance evaluation results show that trajectory similarity can be efficiently supported by these schemes. In all the experiments conducted, the method that uses only one M-tree performs marginally better than the method which utilizes two M-trees. However, the existence of two M-trees offers a higher degree of flexibility during query processing.

Future research may involve: (i) the investigation of alternative indexing schemes, (ii) the study of approximate processing, (iii) the efficient support of trajectory-based *k*-nearest-neighbor processing, and (iv) the utilization of the proposed similarity measures for data mining (e.g., trajectory clustering).

References

- Brinkhoff, T., 2002. A framework for generating network-based moving objects. Geoinformatica 6 (2), 153–180.
- Ciaccia, P., Patella, M., Zezula, P., 1997. M-tree: an efficient access method for similarity search in metric spaces. In: Proceedings of the 23rd International Conference on Very Large Databases (VLDB).
- Faloutsos, C., Ranganathan, M., Manolopoulos, Y., 1994. Fast subsequence matching in time-series databases. In: Proceedings of the ACM SIGMOD Conference.
- Guttman, A., 1984. R-trees: a dynamic index structure for spatial searching. In: Proceedings of the ACM SIGMOD Conference, p. 4757.
- Hwang, J.-R., Kang, H.-Y., Li, K.-J., 2005. Spatio-temporal similarity analysis between trajectories on road networks. In: ER Workshops, pp. 280–289.
- Hwang, J.-R., Kang, H.-Y., Li, K.-J., 2006 Searching for similar trajectories on road networks using spatio-temporal similarity. In: Proceedings of the 10th East European Conference on Advances in Databases and Information Systems (ADBIS), pp. 282–295.
- Jensen, C.S., Kolarvr, J., Pedersen, T.B., Timko, I., 2003. Nearest neighbor queries in road networks. In: Proceedings of the 11th ACM International Symposium on Advances in Geographic Information Systems (ACM GIS).
- Kollios, G., Gounopoulos, D., Tsotras, V.J., 1999a. Nearest neighbor queries in a mobile environment. In: Proceedings of the International Workshop on Spatiotemporal Database Management, pp. 119–134.
- Kollios, G., Gunopoulos, D., Tsotras, V., 1999b. On indexing mobile objects. In: ACM PODS, pp. 261–272.
- Laurinen, P., Siirtola, P., Roning, J., 2006. Efficient algorithm for calculating similarity between trajectories containing an increasing dimension. In: Proceedings of the

24th IASTED International Conference on Artificial Intelligence and Applications, pp. 392–399.

- Lazaridis, I., Porkaew, K., Mehrotra, S., 2002. Dynamic queries over mobile objects. In: EDBT, pp. 269–286.
- Lee, S.-L., Chun, S.-J., Kim, D.-H., Lee, J.-H., Chung, C.-W., 2000. Similarity search for multidimensional data sequences. In: Proceedings of the 16th International Conference on Data Engineering (ICDE).
- Lomet, D., Salsberg, B., 1989. Access methods for multiversion data. In: ACM SIGMOD, pp. 315–324.
- Meratnia, N., de By, R.A., 2002. Aggregation and comparison of trajectories. In: Proceedings of the 10th ACM International Symposium on Advances in Geographic Information Systems (ACM GIS).
- Nascimento, M.A., Silva, J.R.O., 1998. Towards historical R-trees. In: ACM SAC.
- Papadias, D., Zhang, J., Mamoulis, N., 2003. Query processing in spatial network databases. In: Proceedings of the 29th International Conference on Very Large Databases (VLDB).
- Pfoser, D., Jensen, C.S., Theodoridis, Y., 2000. Novel approaches to the indexing of moving object trajectories. In: Proceedings of the 26th International Conference on Very Large Databases (VLDB), pp. 395–406.
- Saltenis, S., Jensen, C.S., Leutenegger, S., Lopez, M., 2000. Indexing the positions of continuously moving objects. In: ACM SIGMOD, pp. 331–342.
- Sankaranarayanan, J., Alborzi, H., Samet, H., 2005. Efficient query processing on spatial networks. In: Proceedings of the 13th ACM International Symposium on Geographic Information Systems (ACM GIS).
- Tao, Y., Papadias, D., 2001a. Efficient historical R-trees. In: Proceedings of the International Conference on Scientific and Statistical Database Management (SSDBM).
- Tao, Y., Papadias, D., 2001b. MV3R-tree a spatio-temporal access method for timestamp and interval queries. In: Proceedings of the 27th International Conference on Very Large Databases (VLDB), pp. 431–440.
- Theodoridis, Y., Sellis, T., Papadopoulos, A.N., Manolopoulos, Y., 1998. Specifications for efficient indexing in spatio-temporal databases. In: Proceedings of the International Conference on Scientific and Statistical Database Management (SSDBM).
- Traina, C., Traina, A.J.M., Seeger, B., Faloutsos, C., 2000. Slim-trees: high performance metric trees minimizing overlap between nodes. In: Proceedings of the Seventh International Conference on Extending Database Technology (EDBT), pp. 51–65.
- Vlachos, M., Gunopulos, D., Kollios, G., 2002a. Robust similarity measures for mobile object trajectories. In: Proceedings of the Fifth International Workshop on Mobility in Databases and Distributed Systems.
- Vlachos, M., Kollios, G., Gunopulos, D., 2002b. Discovering similar multidimensional trajectories. In: Proceedings of the 18th IEEE International Conference on Data Engineering (ICDE).
- Wolfson, O., Xu, B., Chamberlain, S., Jiang, L., 1998. Moving objects databases: issues and solutions. In: Proceedings of the International Conference on Scientific and Statistical Database Management (SSDBM), pp. 111–122.
- Wolfson, O., Xu, B., Chamberlain, S., 2000. Location prediction and queries for tracking moving objects. In: Proceedings of the 16th IEEE International Conference on Data Engineering (ICDE), pp. 687–688.
- Yanagisawa, Y., Akahani, J.-I., Satoh, T., 2003. Shape-based similarity query for trajectory of mobile objects. In: Proceedings of the Fourth International Conference on Mobile Data Management (MDM), pp. 63–77.
- Yoo, J.S., Shekhar, S., 2005. In-route nearest neighbor queries. Geoinformatica 9 (2), 117–137.

Eleftherios Tiakas was born in Thessaloniki, Greece in 1972. He received a B.Sc. in Mathematics (1994), a B.Sc. in Computer Science - Informatics (2006) and a M.Sc. in Information Systems (2006) from the Departments of Mathematics and Informatics of Aristotle University of Thessaloniki, Greece. Currently, he is a Ph.D. candidate at the Department of Informatics of Aristotle University. His research interests include: Spatial & Spatio-Temporal Databases, Information Retrieval, Algorithm Design and Analysis.

Apostolos N. Papadopoulos was born in Eleftheroupolis, Greece in 1971. He received his 5-year Diploma degree in Computer Engineering and Informatics from the University of Patras and his Ph.D. degree from Aristotle University of Thessa-

loniki in 1994 and 2000 respectively. He has published several research papers in journals and proceedings of international conferences. From March 1998 to August 1998 he was a visitor researcher at INRIA Research Center in Rocquencourt, France, to perform research in spatial databases. His research interests include databases, data stream processing, data mining and information retrieval. His research work has over 350 citations in scientific journals and conference proceedings. He has served as a track co-chair of ACM SAC DTTA (Database Technologies Techniques and Applications) Track 2005, 2006, 2007 and 2008. He is a member of the Technical Chamber of Greece. Currently, he is a Lecturer in the Department of Informatics of Aristotle University of Thessaloniki.

Alexandros Nanopoulos was born in Craiova, Romania, in 1974. He graduated from the Department of Informatcis, Aristotle University of Thessaloniki, Greece, on November 1996, and obtained a Ph.D. from the same institute, on February 2003. The subject of his dissertation was: "Techniques for Non Relational Data Mining". He is co-author of more than 30 articles in international journals and conferences. He has also co-author the monograph "Advanced Signature Techniques for Multimedia and Web Applications" and "R-trees: Theory and Applications". His research interests include data mining, web information retrieval, and spatial database indexing.

Yannis Manolopoulos was born in Thessaloniki, Greece in 1957. He received a B.E. (1981) in Electrical Eng. and a Ph.D. (1986) in Computer Eng., both from the Aristotle Univ. of Thessaloniki. Currently, he is Professor at the Department of Informatics of the latter university. He has been with the Department of Computer Science of the University. of Toronto, the Department of Computer Science of the University of Maryland at College Park and the Department of Computer Science of the University of Cyprus. He has published about 200 papers in refereed scientific journals and conference proceedings. He is co-author of the following books: "Advanced Database Indexing" and "Advanced Signature Indexing for Multimedia and Web Applications" by Kluwer, as well as "Nearest Neighbor Search: a Database Perspective" and "R-trees: Theory and Applications" by Springer. His published work has received over 1700 citations from over 450 institutional groups. He served/serves as General/ PC Chair/Cochair of the 8th National Computer Conference (2001), the 6th ADBIS Conference (2002) the 5th WDAS Workshop (2003), the 8th SSTD Symposium (2003), the 1st Balkan Conference in Informatics (2003), the 16th SSDBM Conference (2004) and the 8th ICEIS Conference (2006), the 10th ADBIS Conference (2006). His research interests include Databases. Data mining, Web and Geographical Information Systems, Bibliometrics/Webometrics, Performance evaluation of storage subsystems. Further information can be found at http://delab.csd.auth.gr/manolopo.

Dragan Stojanovic is an Assistant Professor at the Computer Science Department, Faculty of Electronic Engineering, University of Nis, Serbia. He received his Ph.D., M.Sc., and B.Sc. degrees in Computer Science from the University of Nis, in 2004, 1998 and 1993, respectively. His research and development interests encompass context-aware and location-based services, mobile objects and spatio-temporal data management, mobile/Web information systems and services, and geographic information systems. He has published widely in those and related topics. He successfully participates in several international and national R&D projects in cooperation with academic partners and industry.

Slobodanka Dordevic-Kajan is a full professor of computer science and the head of the CG&GIS Lab at the Computer Science Department, Faculty of Electronic Engineering, University of Nis, Serbia. She received her Ph.D., M.S., and B.S. degrees in Computer Science from the Faculty of Electronic Engineering, University of Nis, Serbia, in 1987, 1980 and 1968, respectively. Her current professional and scientific interests include context-aware and location-based services and systems, spatio-temporal and multimedia databases, and application of ontologies to geographic information systems and control and command systems.