

Density-based Place Clustering in Geo-Social Networks*

Jieming Shi, Nikos Mamoulis, Dingming Wu, David W. Cheung
Department of Computer Science, The University of Hong Kong
Pokfulam Road, Hong Kong
{jmshi,nikos,dmwu,dcheung}@cs.hku.hk

ABSTRACT

Spatial clustering deals with the unsupervised grouping of places into clusters and finds important applications in urban planning and marketing. Current spatial clustering models disregard information about the people who are related to the clustered places. In this paper, we show how the density-based clustering paradigm can be extended to apply on places which are visited by users of a geo-social network. Our model considers both spatial information and the social relationships between users who visit the clustered places. After formally defining the model and the distance measure it relies on, we present efficient algorithms for its implementation, based on spatial indexing. We evaluate the effectiveness of our model via a case study on real data; in addition, we design two quantitative measures, called social entropy and community score to evaluate the quality of the discovered clusters. The results show that geo-social clusters have special properties and cannot be found by applying simple spatial clustering approaches. The efficiency of our index-based implementation is also evaluated experimentally.

Categories and Subject Descriptors

H.2.8 [DATABASE MANAGEMENT]: Database Applications—Data mining, Spatial databases and GIS; H.3.3 [INFORMATION STORAGE AND RETRIEVAL]: Information Search and Retrieval—Clustering

Keywords

geo-social network; density-based clustering; spatial indexing

1. INTRODUCTION

Clustering is commonly used as a method for data exploration, characterization, and summarization. Density-based clustering [9], in particular, divides a large collection of points into densely populated regions and it is the most appropriate clustering paradigm for spatial data, which have low dimensionality [30]. Density-based clusters have arbitrary shapes and sizes and exclude objects

in areas of low density (i.e., outliers). The DBSCAN model [9] finds the spatial *eps*-neighborhood of each point p in the dataset, which is a circular region centered at p with radius *eps*. If the *eps*-neighborhood of p is dense, meaning that it contains no less than *MinPts* places, p is called a *core* point. Dense *eps*-neighborhoods are put into the same cluster if they contain the cores of each other.

In this paper, we investigate the extension of traditional density-based clustering for spatial locations to consider their relationship to a social network of people who visit them. In specific, we consider the places of a *Geo-Social Network* (GeoSN) application, which allows users to capture their geographic locations and share them in the social network, by an operation called *checkin*. Online social networks with this functionality include Gowalla¹, Foursquare², and Facebook Places³. A checkin is a triplet $\langle uid, pid, time \rangle$ modeling the fact that user *uid* visited place *pid* at a certain *time*.

We define the new problem of Density-based Clustering Places in Geo-Social Networks (DCPGS), to detect geo-social clusters in GeoSNs. DCPGS extends DBSCAN by replacing the Euclidean distance threshold *eps* for the extents of dense regions by a threshold ϵ , which considers both the *spatial* and the *social* distances between places. For two places p_i and p_j , the spatial distance is considered to be the Euclidean distance between p_i and p_j , while the social distance should consider the social relationships between the two sets of users U_{p_i} and U_{p_j} which have checked in p_i and p_j , respectively. We define and use such a social distance measure, based on the intuition that two places are socially similar if they share many common users in their checkin records or the users in these records are linked by friendship edges. Figure 1(a) illustrates the data of a GeoSN that includes eight users (u_1 – u_8) and two places (p_i and p_j). The dashed lines represent user friendships and the solid lines annotated with timestamps (t_1 – t_9) illustrate checkins of users at places. For instance, user u_1 is a friend with u_2 and u_3 and has visited place p_i at time t_3 and p_j at t_5 . As Figure 1(b) shows, each place is modeled by its spatial coordinates (e.g., $\langle la_i, lo_i \rangle$ for the latitude and longitude of p_i) and the set of users that have visited it (e.g., U_{p_i} for p_i). For the spatial distance between p_i and p_j , we can use the Euclidean distance between $\langle la_i, lo_i \rangle$ and $\langle la_j, lo_j \rangle$, while for the social distance component we use the set of common users in U_{p_i} and U_{p_j} (e.g., $\{u_1, u_2\}$) and the users in one place's record who are friends with visitors of the other place (e.g., $u_3 \in U_{p_i}$ and $u_6 \in U_{p_j}$ who are friends with each other). The intuition is that users who are friends can influence each other to visit the places included in their checkin history. The details of our distance measure are presented in Section 2.

*Work supported by grant HKU 715413E from Hong Kong RGC.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGMOD'14, June 22–27, 2014, Snowbird, UT, USA.
Copyright 2014 ACM 978-1-4503-2376-5/14/06 ...\$15.00
<http://dx.doi.org/10.1145/2588555.2610497>.

¹<http://gowalla.com>

²<https://foursquare.com>

³<https://www.facebook.com/about/location>

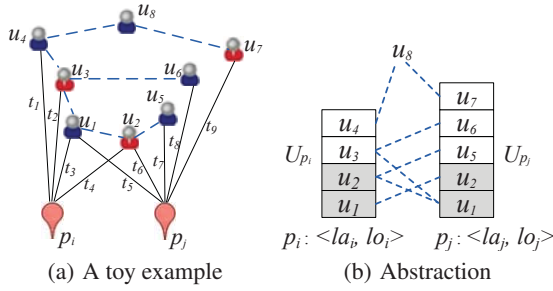


Figure 1: Example and storage structure of GeoSNs

To the best of our knowledge, there is no previous work on clustering GeoSN places. While there exists a significant body of research on analyzing and querying GeoSN data [4, 7, 27, 28, 36, 37], most of these works are centered around users; i.e., they study user behavior, user link prediction or recommendation, or the evaluation of user queries. Thus, the places and checkins are only regarded as some auxiliary information to facilitate user-centered analysis. On the other hand, GeoSNs provide a new and rich form of geographical data, affiliated with the social network graph, the analysis of which can provide new and interesting insights, compared to raw spatial data. In specific, clustering of places in a GeoSN network finds a number of interesting applications:

Generalization and characterization of places. In geographic data analysis, a common task is to define regions (especially in urban areas), which include similar places with respect to the people who live in them or visit them. For example, in urban planning, land managers are interested in identifying regions which have uniform (i.e., consistent) demographic statistics, e.g., areas where elderly people prefer to visit, or people who belong to certain religious communities and have special transportation or living needs. Our DCPGS framework is especially useful for such spatial generalization and characterization tasks, because it can identify geographic regions where places form dense regions and the people who visit them are also socially connected to each other. On the other hand, by simply using spatial clustering, we may not be able to separate regions that are geographically close (e.g., neighboring “districts”) but are visited by different social user groups.

Data cleaning. Intuitively, places that belong in the same geo-social cluster, according to our DCPGS framework, should have similar semantics. Therefore, our clustering results can help toward the cleaning of semantics (e.g., tags), which are given to places being in the same cluster (e.g., inspect tags that are inconsistent with the ones given to the majority of places in the cluster). In addition, as already mentioned, nearby GeoSN locations collected by user checkins could belong to the same physical place (e.g., a large restaurant) and our clusters can help toward identifying such cases and integrating multiple locations to the same physical place (i.e., region), possibly with the joint help of map-matching tools.

Marketing. GeoSN places may be commercial (e.g., restaurants). The fact that two (or more) such places belong to the same geo-social cluster indicates that there is a high likelihood that a user who likes one place would also be interested to visit the other(s). Therefore, by having knowledge of a place’s geo-social cluster, the management of the place may initiate campaigns to users who visited other places in the same cluster, or a set of places could do collaborative promotion (e.g., a discount for users who visit multiple places in the cluster).

Compared to conventional density-based spatial clusters, geo-social clusters detected by DCPGS exhibit larger intra-cluster social strength, as we confirm experimentally in Section 4. DCPGS

also has additional advantages. First, DCPGS uses geo-social splitting criteria; for example, DCPGS splits clusters, which are spatially dense, but they are separated by barriers, such as rivers or walls, or visitors’ weak social connections. Second, DCPGS finds spatially loose clusters that include sets of places that (pairwise) are not very close to each other (and therefore violate a typically tight spatial density threshold), but have very tight social relationships with each other. Such places satisfy our DCPGS criteria. On the other hand, DBSCAN is less flexible in including them in the same cluster as loosening its spatial distance threshold would result in putting everything in a single huge cluster. Third, DCPGS can discover geo-social clusters with fuzzy spatial boundaries; such clusters cannot be identified by spatial clustering, which defines strict spatial boundaries between clusters. Our evaluation is based on case studies and on the use of quantitative measures that we also propose in this paper. In addition, we demonstrate that the social distance measure we propose and use in DCPGS is more effective compared to alternative measures (based on node-to-node proximity). Overall, the results of our evaluation indicate that the social relationships between users who visit places have great impact in the clustering of places and cannot be overlooked.

In addition to the effectiveness of place clustering according to our DCPGS framework, we are also concerned about the efficient implementation of DCPGS in this paper. Real GeoSNs generate millions of checkins and contain millions of places, the efficient clustering of which is challenging. We present an intuitive implementation of DCPGS, which applies repetitive range searches on a traditional spatial index (i.e., an R-tree), extending the original implementation of DBSCAN. In view of its high cost, we propose a much more efficient implementation, which relies on a dynamic grid partitioning technique, used to efficiently compute densely connected spatial neighborhoods of places. This grid-based implementation can cluster millions of GeoSN places in just a few seconds. Besides, the experiments demonstrate that our proposed social distance measure between places is very efficient to compute.

Summing up, the contributions of this paper are as follows:

- We propose and formulate the problem of density-based clustering GeoSN places.
- We define a simple but effective social distance measure between places in GeoSNs.
- We design efficient algorithms for implementing DCPGS.
- We demonstrate the effectiveness of DCPGS by case studies and quantitative evaluation through two quality measures that are also devised in this paper.
- We experimentally analyze the efficiency of our methods.

The rest of the paper is organized as follows. Section 2 formulates the DCPGS problem and defines the social distance measure between places that we use. DCPGS algorithms based on R-tree and grid partitioning are proposed in Section 3. The effectiveness and efficiency of our framework are analyzed in Sections 4 and 5, respectively. Related work is reviewed in Section 6. Finally, Section 7 concludes this paper and discusses future work.

2. MODEL AND DEFINITIONS

Our data input includes three components: a *social network*, a set of *places* and the *checkins* of users to the places. The social network is an undirected graph $G = (U, E)$, where U is the set of all users and each edge $(u_i, u_j) \in E$ indicates that users $u_i, u_j \in U$ are friends. Set P contains the set of all places visited by users, in the form of $\langle \text{latitude}, \text{longitude} \rangle$ GPS points. Thus, identifiers are assigned to places according to their distinct GPS coordinates. Set $CK = \{(u_i, p_k, t_r) | u_i \in U \text{ and } p_k \in P\}$ includes all checkins generated

by users in U . For a place p_k , the set U_{p_k} of *visiting users* of p_k is defined by $U_{p_k} = \{u_i | \langle u_i, p_k, * \rangle \in CK\}$, where $*$ means any time. Figure 1(b) shows U_{p_i} and U_{p_j} for the two places p_i and p_j of the toy example in Figure 1(a). The figure also connects the user pairs in the two sets who are linked by friendship edges in the social network. Note that user u_8 does not belong to either U_{p_i} or U_{p_j} , but connects users u_4 and u_7 in the social graph.

2.1 DCPGS Model

Our Density-based Clustering Places in Geo-Social Networks (DCPGS) model extends the model of DBSCAN [9]; for each place p_i in the GeoSN, DCPGS finds the *geo-social ϵ -neighborhood* $N_\epsilon(p_i)$ of p_i , which includes all places p_j such that $D_{gs}(p_i, p_j) \leq \epsilon$, $D_S(p_i, p_j) \leq \tau$, and $E(p_i, p_j) \leq \max D$. For two places p_i, p_j , $E(p_i, p_j)$ is the *Euclidean distance*, $D_S(p_i, p_j)$ is the *social distance*, and $D_{gs}(p_i, p_j) = f(D_S(p_i, p_j), E(p_i, p_j))$ is the *geo-social distance*, defined as a function of $E(p_i, p_j)$ and $D_S(p_i, p_j)$. Parameter ϵ is geo-social distance threshold, while τ and $\max D$ are two *sanity* constraints for the social and the spatial distances between places, respectively. We will give detailed definitions for all above distance functions and parameters later on. If the geo-social ϵ -neighborhood of a place p_i contains at least $MinPts$ places, then p_i is a *core place*; in this case, p_i and all places in its geo-social ϵ -neighborhood should belong to a cluster $r(p_i)$. If another core place p_j belongs to cluster $r(p_i)$, then $r(p_i) = r(p_j)$, i.e., the clusters defined by p_i and p_j are merged. After identifying all core places and merging the corresponding clusters, DCPGS ends up with a set of (disjoint) clusters and a set of outliers (i.e., places that do not belong to the geo-social ϵ -neighborhood of any core place). In Section 3, we present algorithms for generating DCPGS clusters.

Parameters. ϵ and $MinPts$ are the main parameters of DCPGS. $MinPts$ (i.e., the minimum number of places in the neighborhood of a core point) is set as in the original DBSCAN model (see [9]); a typical value is 5. ϵ takes a value between 0 and 1, because, as we explain later on, we define $D_{gs}(p_i, p_j)$ to take values in this range. Since the geo-social distance $D_{gs}(p_i, p_j)$ is a function of a spatial and a social distance, τ and $\max D$ constrain these individual distances to avoid the following two cases that negatively affect the quality of geo-social clusters.

- The geo-social distance between two places p_i and p_j could be less than ϵ if they are extremely close to each other in space, but have no social connection at all. This may lead to putting places close to each other spatially, but having no social relationship, into the same cluster.
- The geo-social distance between two places p_i and p_j could be less than ϵ if they have very small social distance, but they are extremely far from each other spatially. This may lead to putting places with close social distances, but large spatial distances, into the same cluster.

Constraints τ and $\max D$ are defined for quality control and can be set by experts or according to the analyst's experience. We experimentally study how clustering quality is affected by the two constraints and ϵ in Section 4.

Distance Functions. The social distance $D_S(p_i, p_j)$ takes as inputs the sets of users U_{p_i} and U_{p_j} who have visited p_i and p_j , respectively, and returns a value between 0 and 1. In Section 2.2, we present our definition for $D_S(p_i, p_j)$ and alternative ways to define it based on previous work. Before defining the geo-social distance $D_{gs}(p_i, p_j)$, we *normalize* the Euclidean distance $E(p_i, p_j)$ to a *spatial distance* $D_P(p_i, p_j) = \frac{E(p_i, p_j)}{\max D}$ that takes values between 0 and 1. Finally, $D_{gs}(p_i, p_j)$ is defined as weighted sum of

$D_S(p_i, p_j)$ and $D_P(p_i, p_j)$, i.e.,

$$D_{gs}(p_i, p_j) = \omega \cdot D_P(p_i, p_j) + (1 - \omega) \cdot D_S(p_i, p_j), \quad (1)$$

where $\omega \in [0, 1]$.

2.1.1 Alternatives to DCPGS

Our place clustering model (DCPGS) extends density-based clustering in spatial databases. GeoSN places can alternatively be clustered by the use of graph clustering models. The main idea of such a model is to construct a *place network* PN , which connects places according to their social and spatial distances and then apply an off-the-shelf community detection algorithm on PN . Specifically, given two places p_i and p_j , if $E(p_i, p_j) \leq \max D$ and $D_S(p_i, p_j) \leq \tau$, an *undirected* and *weighted* edge with geo-social weight $W_{gs}(p_i, p_j) = 1 - D_{gs}(p_i, p_j)$ is added between p_i and p_j . Community detection algorithms like Link Clustering [1, 6] or Metis [16] can then be applied to derive the clusters. Link Clustering constructs a dendrogram of network communities (that may overlap) in a hierarchical manner. Metis is another multilevel graph partition paradigm that includes three phases: graph coarsening, initial partitioning, and uncoarsening; Metis divides a network into k non-overlapping communities. As we will show in Section 4, these graph clustering methods are inferior to DCPGS.

2.2 Social Distance Between Places

The social distance $D_S(p_i, p_j)$ between p_i and p_j naturally depends on the social network relationships between the sets U_{p_i} and U_{p_j} of users who visited p_i and p_j , respectively. Our definition for $D_S(p_i, p_j)$ is based on the set CU_{ij} of *contributing users* between two places p_i and p_j :

Definition 1. (Contributing Users) Given two places p_i and p_j with visiting users U_{p_i} and U_{p_j} , respectively, the set of contributing users CU_{ij} for the place pair (p_i, p_j) is defined as

$$CU_{ij} = \{u_a \in U_{p_i} | u_a \in U_{p_j} \text{ or } \exists u_b \in U_{p_j}, (u_a, u_b) \in E\} \\ \cup \{u_a \in U_{p_j} | u_a \in U_{p_i} \text{ or } \exists u_b \in U_{p_i}, (u_a, u_b) \in E\} \quad (2)$$

Specifically, if a user u_a has visited both p_i and p_j , then u_a is a contributing user. Also if u_a has visited place p_i , u_b has visited p_j , and u_a and u_b are friends, both u_a and u_b are contributing users. Users in CU_{ij} contribute positively (negatively) to the social similarity (distance) between p_i and p_j . Formally:

Definition 2. (Social Distance) Given two places p_i and p_j with visiting users U_{p_i} and U_{p_j} , respectively, the *social distance* between p_i and p_j is defined as

$$D_S(p_i, p_j) = 1 - \frac{|CU_{ij}|}{|U_{p_i} \cup U_{p_j}|} \quad (3)$$

The above definition of $D_S(p_i, p_j)$ takes both the set similarity between sets U_{p_i} and U_{p_j} and the social relationships among users in U_{p_i} and U_{p_j} into account. In addition, the distance measure penalizes pairs of places p_i and p_j which are popular (i.e., U_{p_i} and/or U_{p_j} are large) but their set of contributing users is relatively small (see Equation 3). The reason is that such place pairs are not characteristic to their (loose) social connections. As an example, consider places p_i and p_j of Figure 1. To compute $D_S(p_i, p_j)$, we first set $U_{p_i} = \{u_1, u_2, u_3, u_4\}$ and $U_{p_j} = \{u_1, u_2, u_5, u_6, u_7\}$. All users in U_{p_i} and U_{p_j} are checked one by one to obtain the contributing users between p_i and p_j . We derive $CU_{ij} = \{u_1, u_2, u_3, u_5, u_6\}$, since (i) both u_1 and u_2 have visited p_i and p_j , (ii) user u_3 , who visited p_i , has a friend u_6 who visited p_j , (iii) symmetrically, user

u_6 , who visited p_j , has a friend u_3 who visited p_i , and (iv) u_5 ($\in U_{p_j}$) has a friend u_2 having been to p_i . According to Definition 2, the social distance $D_S(p_i, p_j)$ between p_i and p_j in Figure 1 is $1 - |CU_{ij}|/(|U_{p_i} \cup U_{p_j}|) = 1 - 5/7 \approx 0.2857$.

Observe that only direct friendship edges between users of U_{p_i} and U_{p_j} are considered in our social distance definition. Longer network paths, such as friend-of-friend relationships, are ignored (e.g., the case of users u_4 and u_7 in Figure 1(b) who are connected via user u_8). According to the small world effect [20], a user in a social network can reach a large portion of other users within only few hops. For instance, the 90-percentile effective diameter [17] of Gowalla GeoSN, used in our experiments, is just 5.7.⁴ This means within quite a few hops most users can reach a very large percentage of all users. In Gowalla, only 8 users can access more than 1% of all the users in 1 hop, while 40516 users (20.61% of all the users) can access more than 1% of all the users in 2 hops and 141582 users (72.02% of all the users) can reach more than 1% of all the users in 3 hops. The number of users who can reach more than 1% users in 2 or 3 hops increases dramatically compared to the percentage of those visited within 1 hop. Thus, paths longer than 1 hops are too common and cannot be considered as (indirect) user relationships; i.e., their impact is much weaker compared to direct friendship edges. Hence, Definition 2 introduces a simple, but powerful social distance measure. Properties of $D_S(p_i, p_j)$ include *symmetry* (i.e., $D_S(p_i, p_j) = D_S(p_j, p_i)$) and *self-minimality* (i.e., $D_S(p_i, p_j) \in [0, 1]$, and $D_S(p_i, p_j) = 0$ for $U_{p_i} = U_{p_j}$). On the other hand, $D_S(p_i, p_j)$ does not obey the triangular inequality, but this does not affect our clustering algorithm.

2.2.1 Alternatives to D_S

Our DCPGS model is independent of the social distance definition between places (i.e., D_S). As an alternative to our Definition 2, the following measures can be used. In Section 4, we evaluate the effectiveness of these alternatives.

Jaccard. Based on the Jaccard similarity $J(p_i, p_j) = (|U_{p_i} \cap U_{p_j}|)/(|U_{p_i} \cup U_{p_j}|)$ between the sets of visiting users for p_i and p_j , we can define the following distance:

$$D_S^{Jacc}(p_i, p_j) = 1 - J(p_i, p_j)$$

$D_S^{Jacc}(p_i, p_j)$ is not intuitive; it disregards the social network, assuming that two users who are friends do not affect each other in visiting GeoSN places.

SimRank. SimRank is a structural-context model for measuring the similarity between nodes in a graph. The idea is that two nodes are equivalent if they relate to equivalent nodes. We can define a SimRank-based social distance $D_S^{sim}(p_i, p_j)$, using the Minimax version of SimRank [14].⁵ This measure compares each of p_i 's visiting users $u_r^{p_i}$ with the visiting user $u_s^{p_j}$ of p_j who is the most similar to $u_r^{p_i}$, to compute the similarity between places $s(p_i, p_j)$. The similarity between users $s(u_r^{p_i}, u_s^{p_j})$ is computed in an analogous way. Specifically,

$$D_S^{sim}(p_i, p_j) = 1 - s(p_i, p_j),$$

where $s(p_i, p_j) = \min(s_{p_i}(p_i, p_j), s_{p_j}(p_i, p_j))$,

$$s_{p_i}(p_i, p_j) = \frac{\phi}{|U_{p_i}|} \sum_{u_r \in U_{p_i}} \max_{u_s \in U_{p_j}} s(u_r, u_s),$$

where $\phi = 0.8$ is a decay factor [14],

$s(u_r, u_s) = \min(s_{u_r}(u_r, u_s), s_{u_s}(u_r, u_s))$, and assuming that

⁴<http://snap.stanford.edu/data/index.html>

⁵The original SimRank measure is only meant for node-to-node similarity; in our case, we need a measure between U_{p_i} and U_{p_j} .

P_{u_r} is the set of places visited by u_r ,

$$s_{u_r}(u_r, u_s) = \frac{\phi}{|P_{u_r}|} \sum_{p_i \in P_{u_r}} \max_{p_j \in P_{u_s}} s(p_i, p_j).$$

Katz. The Katz similarity measure [34] sums over all possible paths from user u_r to u_s with exponential damping by length, i.e.,

$$\mathcal{K}(u_r, u_s) = \sum_{l=1}^{\infty} \beta^l |paths_{u_r, u_s}^l|, \text{ where } paths_{u_r, u_s}^l \text{ is the set of all}$$

length- l paths from u_r to u_s , and damping factor β is typically set to 0.05. Due to the poor scalability of this measurement, in practice only paths up to length L are considered [33]; i.e., an approximated

$$\text{Katz score } \mathcal{K}_a(u_r, u_s) = \sum_{l=1}^L \beta^l |paths_{u_r, u_s}^l| \text{ can be used. Accord-}$$

ingly, we can define a Katz-based social distance between places p_i and p_j , by averaging the normalized Katz similarities between all pairs of users from U_{p_i} and U_{p_j} :

$$D_S^{Katz}(p_i, p_j) = 1 - \frac{1}{|U_{p_i}| |U_{p_j}|} \sum_{u_r \in U_{p_i}} \sum_{u_s \in U_{p_j}} \mathcal{K}_a(u_r, u_s)$$

CommuteTime. The hitting time $h(u_r, u_s)$ from u_r to u_s is the expected number of steps required for a random walk starting at u_r to reach u_s . The commute time between u_r and u_s is defined by $ct(u_r, u_s) = h(u_r, u_s) + h(u_s, u_r)$. However, the commute time is sensitive to long paths and favors nodes of high degree. Thus, the *truncated commute time* [26], which considers only paths of length no longer than L , can be used to model the social distance between a pair of users. Finally, we can define a commute time based social distance between places p_i and p_j as follows:

$$D_S^{ct}(p_i, p_j) = \frac{1}{|U_{p_i}| |U_{p_j}|} \sum_{u_r \in U_{p_i}} \sum_{u_s \in U_{p_j}} ct^L(u_r, u_s)$$

where $ct^L(u_r, u_s)$ is the normalized truncated commute time.

3. ALGORITHMS

We propose two algorithms for DCPGS. DCPGS-R (Section 3.1) is based on the R-tree index, while DCPGS-G (Section 3.2) uses a grid partitioning.

3.1 Algorithm DCPGS-R: R-tree based

Algorithm DCPGS-R is a direct extension of the DBSCAN algorithm; it uses an R-tree to facilitate the search of the geo-social ϵ -neighborhood for a given place. Initially, all places in the GeoSN are bulk-loaded into an R-tree. Then, DCPGS-R examines all places and, given a place p_i , it performs a range query centered at p_i with radius $maxD$ to get a set of *candidate places* that may fall in the geo-social ϵ -neighborhood of p_i , i.e., $N_\epsilon(p_i)$. Recall that $maxD$ is the maximum allowed spatial distance between place p_i and places in its geo-social ϵ -neighborhood. Then, DCPGS-R keeps in $N_\epsilon(p_i)$ only the candidates that satisfy the social distance constraint τ and the geo-social distance threshold ϵ .

For the sake of efficiency, the social network is stored in a hash table. Specifically, each pair of friends in the social network is recorded as an entry in the hash table, such that checking whether two users are friends or not only incurs constant cost. In addition, for each place p_i , we keep track of its visiting users U_{p_i} . The computation of the social distance (Definition 2) between two places p_i and p_j involves finding the pairs of friends between sets U_{p_i} and U_{p_j} and has insignificant cost compared to the range queries used to compute the set of candidate places.

Algorithm 1 is the pseudocode of DCPGS-R. The identity of the current cluster cid is initialized to 1 in line 1. Queue Q (initialized in line 2) stores the places that have the potential to be added to the

current cluster. Hash table H records whether the geo-social distances between pairs of places have been computed before (line 3) and its use will be explained later. For each *unprocessed* place p_i , its geo-social ϵ -neighborhood $N_\epsilon(p_i)$ is obtained by calling function $\text{GETNEIGH}(p_i, \epsilon, \tau, \text{maxD}, \text{MinPts}, \omega, H)$, outlined in Algorithm 2. A place is *unprocessed* if its geo-social ϵ -neighborhood has not been computed before. If $N_\epsilon(p_i)$ contains at least MinPts places, then p_i is a *core* place, i.e., p_i belongs to a cluster and all places in $N_\epsilon(p_i)$ should be given the same cid as p_i (lines 6-11). Next, all unprocessed places in $N_\epsilon(p_i)$ are pushed into Q for later processing. Lines 12-20 expand the current cluster cid as much as possible by checking the geo-social ϵ -neighborhood of the unprocessed places in Q . No more places can be included in the current cluster when Q is empty. In this case, the algorithm proceeds to find the next cluster (cid is increased in line 21).

Algorithm 1 DCPGS-R(GeoSN, $\epsilon, \tau, \text{maxD}, \text{MinPts}, \omega$)

```

1:  $\text{cid} = 1$ 
2:  $Q = \text{empty}$ 
3: Geo-social distance cache  $H$ 
4: for each unprocessed place  $p_i$  in GeoSN do
5:    $N_\epsilon(p_i) = \text{GETNEIGH}(p_i, \epsilon, \tau, \text{maxD}, \text{MinPts}, \omega, H)$ 
6:   if  $|N_\epsilon(p_i)| \geq \text{MinPts}$  then
7:     assign  $\text{cid}$  to  $p_i$ 
8:     for each place  $p_j \in N_\epsilon(p_i)$  do
9:       assign  $\text{cid}$  to  $p_j$ 
10:      if  $p_j$  is unprocessed then
11:         $Q.\text{push}(p_j)$ 
12:   while  $!Q.\text{isEmpty}()$  do
13:      $p_k = Q.\text{pop}()$ 
14:     if  $p_k$  is unprocessed then
15:        $N_\epsilon(p_k) = \text{GETNEIGH}(p_k, \epsilon, \tau, \text{maxD}, \text{MinPts}, \omega, H)$ 
16:       if  $|N_\epsilon(p_k)| \geq \text{MinPts}$  then
17:         for each place  $p_m \in N_\epsilon(p_k)$  do
18:           assign  $\text{cid}$  to  $p_m$ 
19:           if  $p_m$  is unprocessed then
20:              $Q.\text{push}(p_m)$ 
21:    $\text{cid} = \text{cid} + 1$ 

```

Function GETNEIGH , shown in Algorithm 2, is used to get the geo-social ϵ -neighborhood of place p_i . Initially $N_\epsilon(p_i)$ is empty. Lines 2-4 first perform a spatial range query centered at the current place p_i with radius maxD to get a candidate set CandSet containing the places that may fall in $N_\epsilon(p_i)$. If the size of CandSet is less than MinPts , $N_\epsilon(p_i)$ definitely includes less than MinPts places and, therefore, p_i is a non-core place. Otherwise, the algorithm tries to compute the social distance between every point $p_j \in \text{CandSet}$ and p_i . However, before this, given the fact that the spatial distance between p_i and every candidate place p_j is already obtained in the spatial range query step, a *spatial filter* is employed to avoid unnecessary social distance computations (line 6).

PROPOSITION 1. Spatial filter. *Given two places p_i and p_j with spatial distance $D_P(p_i, p_j)$, if $\omega \cdot D_P(p_i, p_j) > \epsilon$, then $\omega \cdot D_P(p_i, p_j) + (1 - \omega) \cdot D_S(p_i, p_j) > \epsilon$.*

PROOF. Since $\omega \in [0, 1]$ and $D_S(p_i, p_j) \in [0, 1]$, $(1 - \omega) \cdot D_S(p_i, p_j) \geq 0$. Consequently, if $\omega \cdot D_P(p_i, p_j) > \epsilon$, then $\omega \cdot D_P(p_i, p_j) + (1 - \omega) \cdot D_S(p_i, p_j) > \epsilon$. \square

Note that distances D_{gs} and D_S and D_P are all symmetric. Hence, given a place p_i , if p_j is (not) in $N_\epsilon(p_i)$, then p_i is also (not) in $N_\epsilon(p_j)$, and vice versa. Therefore, we keep track in a hash table H (line 14) whether p_i and p_j are in each other's geo-social ϵ -neighborhood once their geo-social distance has been computed. This information is used when computing the geo-social ϵ -neighborhood of p_j later (line 7-9), in order to avoid computing

the distance between the same pair of places twice. Lines 11-15 compute distances, verify candidates, and update H if necessary. Line 16 is another filter, called *sufficient filter*, to check whether there are enough remaining candidate places in CandSet to render p_i 's geo-social ϵ -neighborhood dense. If no, the algorithm can stop without checking the remaining candidate places and return.

Algorithm 2 $\text{GETNEIGH}(p_i, \epsilon, \tau, \text{maxD}, \text{MinPts}, \omega, H)$

```

1:  $N_\epsilon(p_i) \leftarrow \emptyset$ 
2:  $\text{CandSet} = \text{RANGEQUERY}(p_i, \text{maxD})$ 
3: if  $|\text{CandSet}| < \text{MinPts}$  then
4:   return  $N_\epsilon(p_i)$ 
5: for each place  $p_j \in \text{CandSet}$  do
6:   if  $\omega \cdot D_P(p_i, p_j) \leq \epsilon$  then
7:     if  $H.\text{exists}((p_i, p_j))$  then
8:       if  $H[(p_i, p_j)]$  is TRUE then
9:          $N_\epsilon(p_i).\text{insert}(p_j)$ 
10:    else
11:      Compute  $D_S(p_i, p_j)$  and  $D_{gs}(p_i, p_j)$ 
12:      if  $D_S(p_i, p_j) \leq \tau$  &&  $D_{gs}(p_i, p_j) \leq \epsilon$  then
13:         $N_\epsilon(p_i).\text{insert}(p_j)$ 
14:         $H[(p_i, p_j)] \leftarrow \text{TRUE}$ 
15:    $\text{CandSet}.\text{erase}(p_j)$ 
16:   if  $|\text{CandSet}| + |N_\epsilon(p_i)| < \text{MinPts}$  then break
17: return  $N_\epsilon(p_i)$ 

```

3.2 Algorithm DCPGS-G: Grid based

DCPGS-R conducts a spatial range query for each place in the GeoSN to obtain the candidate places for the purpose of discovering geo-social clusters. Even though individual R-tree based range queries are very efficient, discovering geo-social clusters in a GeoSN with millions of places requires millions of such queries (e.g., there are 1,280,969 places in the Gowalla dataset used in our experiments). Given two places p_i and p_j that are spatially close to each other, as Figure 2(a) shows, the results of the two range queries with radius maxD centered at p_i and p_j , respectively, are almost identical. In algorithm DCPGS-R, to get the candidate places CandSet for each place in Figure 2(a), 8 independent range queries are issued on the R-tree that search almost the same space, resulting in redundant traversing paths and computations. To overcome this drawback, we develop a dynamic grid partitioning technique and a new algorithm DCPGS-G.

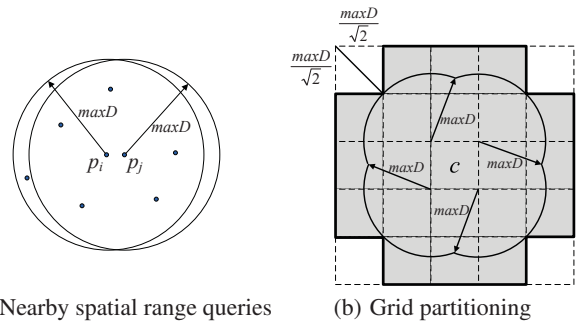


Figure 2: Nearby spatial range queries and grid partitioning

Grid Partitioning. The area covered by the dataset is partitioned by a regular grid with cells of size $\text{maxD}/\sqrt{2} \times \text{maxD}/\sqrt{2}$. The non-empty cells of the grid are indexed by a hash table with the grid cell coordinates as search keys.

Neighbor Cells. The neighbor cells of a cell c are the cells that intersect the union of four circles, each centered at a corner of cell

c with radius $maxD$. For example, in Figure 2(b), the 20 gray cells (except c) are the neighbor cells of c , denoted as $NC(c)$. We can trivially show that for any place p inside c , the content of p 's geo-social ϵ -neighborhood is contained in $NC(c)$ and c itself.

Cluster Discovery. Algorithm 3 is a pseudocode for DCPGS-G. It includes three phases. First, the algorithm maps all places into grid cells (line 1). Second, it obtains the geo-social ϵ -neighborhoods of all places (lines 2-6). The third phase discovers all geo-social clusters in the GeoSN (line 7). The geo-social ϵ -neighborhoods of places are computed at the grid cell level. Specifically, for each nonempty and *unprocessed* cell c , function GETNEIGHCELLS retrieves its neighbor cells $NC(c)$ (GETNEIGHCELLS is trivial and thus details are omitted). A cell is 'unprocessed' if its neighbor cells have not been retrieved before. Function COMPCELLPAIR (Algorithm 4) first filters out the pairs of places (p_i, p_j) with spatial distance greater than $maxD$, where $p_i \in c, p_j \in NC(c)$ and $p_i \neq p_j$. This step is not needed if p_i and p_j are in the same cell (in this case, the spatial distance between p_i and p_j is certainly at most $maxD$). Next, the pairs of places (p_i, p_j) that satisfy the social distance constraint τ and the geo-social distance threshold ϵ are selected. If p_i and p_j are in each others' geo-social ϵ -neighborhood, their corresponding $N_\epsilon(p_i)$ and $N_\epsilon(p_j)$ are updated. After all cells have been processed, meaning that the geo-social ϵ -neighborhoods of all places in the GeoSN are acquired, function GETCLUSTERS is called to discover geo-social clusters following the framework of DCPGS-R (Algorithm 1), except that the $N_\epsilon(p_i)$ of each place p_i has already been computed. Note that an *unprocessed* place p_i in function GETCLUSTERS means that the size of p_i 's geo-social ϵ -neighborhood, $|N_\epsilon(p_i)|$, has not been checked.

Algorithm 3 DCPGS-G(GeoSN, $\epsilon, \tau, maxD, MinPts, \omega$)

```

1: Map all places into grid cells
2: for each nonempty && unprocessed cell  $c$  do
3:   COMPCELLPAIR( $c, c, \epsilon, \tau, maxD, \omega$ )
4:    $NC(c)$  = GETNEIGHCELLS( $c$ )
5:   for each nonempty && unprocessed cell  $c' \in NC(c)$  do
6:     COMPCELLPAIR( $c, c', \epsilon, \tau, maxD, \omega$ )
7: GETCLUSTERS( $N_\epsilon, MinPts$ )

```

Complexity. With the help of grid partitioning, the geo-social ϵ -neighborhood of all places in cell c can be obtained by checking all c 's neighbor cells; the whole process can be completed within a single pass of the data. Thus, the complexity of DCPGS-G is $O(n)$, as each of its three phases makes one pass over the data. However, algorithm DCPGS-R computes the geo-social ϵ -neighborhoods of each place one by one. Hence its cost is $O(n \log n)$, given that the expected cost of a single range query on the R-tree is $O(\log n)$.

Algorithm 4 COMPCELLPAIR(cell c , cell $c', \epsilon, \tau, maxD, \omega$)

```

1: for each pair  $(p_i, p_j)$  where  $p_i \in c, p_j \in c', p_i \neq p_j$  do
2:   if  $c = c' \wedge E(p_i, p_j) \leq maxD$  then
3:     Compute  $D_P(p_i, p_j)$ 
4:     if  $\omega \cdot D_P(p_i, p_j) \leq \epsilon$  then
5:       Compute  $D_S(p_i, p_j)$  and  $D_{gs}(p_i, p_j)$ 
6:       if  $D_S(p_i, p_j) \leq \tau$  &&  $D_{gs}(p_i, p_j) \leq \epsilon$  then
7:          $N_\epsilon(p_i).insert(p_j)$ 
8:          $N_\epsilon(p_j).insert(p_i)$ 

```

4. QUALITATIVE ANALYSIS

This section analyzes the quality of the geo-social clusters discovered by our proposed DCPGS framework. First, we compare

Algorithm 5 GETCLUSTERS($N_\epsilon, MinPts$)

```

1:  $cid = 1$ 
2:  $Q = empty$ 
3: for each unprocessed place  $p_i$  in GeoSN do
4:   if  $|N_\epsilon(p_i)| \geq MinPts$  then
5:     assign  $cid$  to  $p_i$ 
6:     for each place  $p_j \in N_\epsilon(p_i)$  do
7:       assign  $cid$  to  $p_j$ 
8:       if  $p_j$  is unprocessed then
9:          $Q.push(p_j)$ 
10:  while  $!Q.isEmpty()$  do
11:     $p_k = Q.pop()$ 
12:    if  $p_k$  is unprocessed then
13:      if  $|N_\epsilon(p_k)| \geq MinPts$  then
14:        for each place  $p_m \in N_\epsilon(p_k)$  do
15:          assign  $cid$  to  $p_m$ 
16:          if  $p_m$  is unprocessed then
17:             $Q.push(p_m)$ 
18:     $cid = cid + 1$ 

```

DCPGS with the graph clustering approaches discussed in Section 2.1.1, in order to demonstrate the suitability of the density-based clustering model for this application. Second, we compare with two extreme versions of DCPGS: PureSocialDistance applies density-based clustering by using the social distance $D_S(p_i, p_j)$ only, while DBSCAN uses only the Euclidean distance $E(p_i, p_j)$. This comparison shows the appropriateness of using both social and spatial distances in clustering. In the implementation of PureSocialDistance, we do not put place pairs with spatial distance more than 1000m in the same cluster; otherwise this method becomes too expensive. Finally, we assess the suitability of our social distance measure (Section 2.2) by evaluating versions of DCPGS, which use the alternative social distance definitions discussed in Section 2.2.1. All tested methods were implemented in C++ and the experiments were performed on a 3.4 GHz quad-core machine running Ubuntu 12.04 with 16 GBytes memory.

Data. We use two publicly available datasets⁶ from historical geo-social networks. Gowalla contains a social network with $|U| = 196,591$ users and $|E| = 950,327$ undirected friendship edges. There are $|CK| = 6,442,892$ checkins performed by those users on $|P| = 1,280,969$ places over a period from Feb. 2009 to Oct. 2010. Brightkite includes a social network of $|U| = 58,228$ users and $|E| = 214,078$ undirected friendship edges. It contains $|CK| = 4,491,143$ checkins on $|P| = 772,783$ distinct places collected over the period from Apr. 2008 to Oct. 2010.

Default Parameter Settings. The density requirement of the clustering is determined by parameters $MinPts$ and ϵ (or DBSCAN's eps). We set $MinPts = 5$ for all approaches; various density settings can be achieved by just tuning ϵ (or DBSCAN's eps). For instance, a large $MinPts$ has similar effect as a small ϵ (or DBSCAN's eps). By default, parameter ω in the geo-social distance is set to 0.5 to equally weigh the social and spatial distances. By default, parameter τ is set to 0.7, and $maxD$ is set to 100m for dataset Gowalla and 120m for dataset Brightkite.

4.1 Visualization-based Analysis

We first visualize and compare the clusters found by DCPGS and alternative approaches in the area of Manhattan, on the Gowalla dataset. Figures 3(a)-(c) show the clusters by DCPGS, DBSCAN (which disregards the social network behind the places) and PureSocialDistance (which disregards the spatial information). DCPGS finds geo-social clusters with the following features.

⁶downloaded from snap.stanford.edu/data/index.html

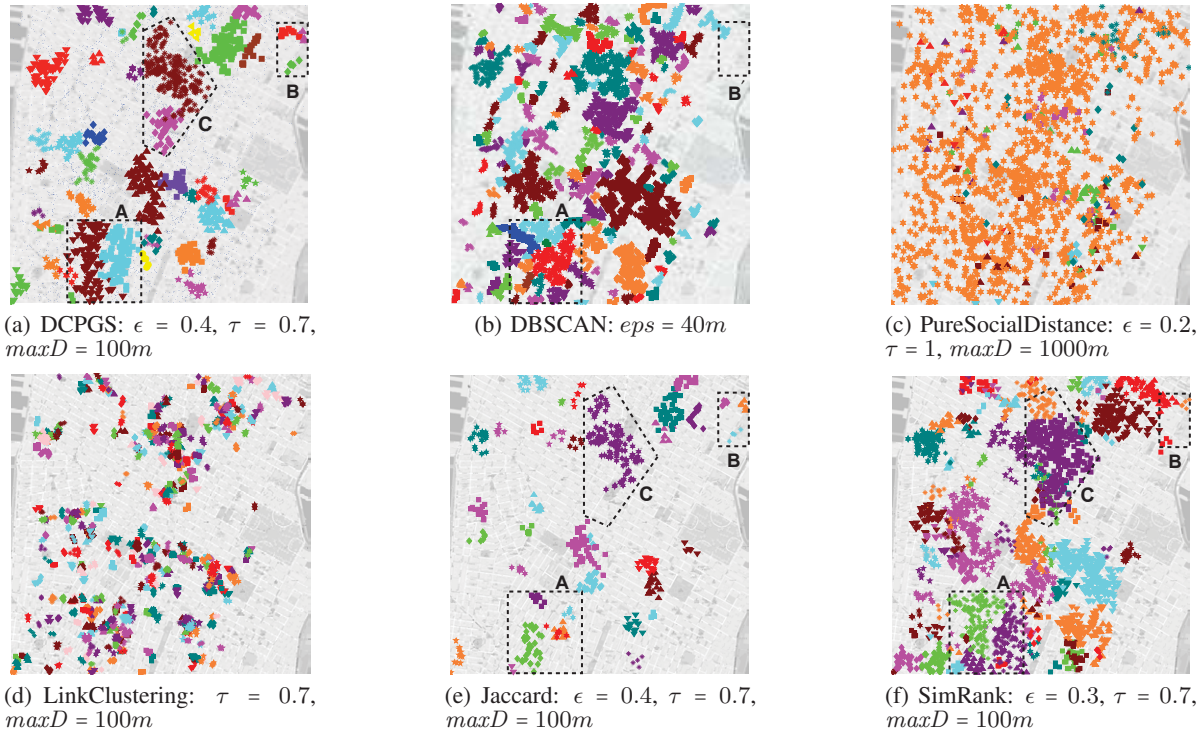


Figure 3: Place clusters of Gowalla found in Manhattan

Geo-Social Splitting/Merging Criteria. Geo-social clusters that are very close to each other are split correctly by DCPGS, while DBSCAN may consider them as a single cluster due to their spatial closeness; in other cases, clusters split by DBSCAN due to relatively low spatial density between them are merged by DCPGS because of their strong social ties. For example, consider region A in Figures 3(a) and the corresponding region A' in Figure 3(b), where DCPGS and DBSCAN detect clusters with totally different layouts. By tuning the parameters of DBSCAN, we are not able to find the clusters found by DCPGS, because the densities of the two clusters in region A are similar and the two clusters are close to each other. Thus, DBSCAN can only consider the places in region A' as either a single cluster or as several fragmented clusters (Figure 3(b)), under different parameter settings. In certain cases, spatially dense clusters may be split by DCPGS because of some natural barriers, such as rivers, and walls. These barriers make it inconvenient to travel from one side to the other, resulting in a splitting effect. As an example, in Figure 4, a cluster (region D) found by DBSCAN is split into two DCPGS clusters (regions D₁ and D₂) by the river, since the users on different river sides are proved to have weak social connection. While it is possible for DBSCAN to find the two DCPGS clusters by reducing the value of ϵ_{ps} , its parameter settings in this case make some existing clusters disappear, resulting in too many outliers.

Spatially Loose Clusters. Some geo-social clusters detected by DCPGS in region B of Figure 3(a) are considered as outliers by DBSCAN in the corresponding region B' of Figure 3(b). Region B' is spatially too sparse to satisfy the density requirement of DBSCAN, and thus most places inside it are filtered out as outliers. However, the users who checked in those places have strong social relationships. Hence, geo-social clusters are discovered in region B by DCPGS in Figure 3(a). While it is possible for DBSCAN to discover such spatially loose clusters by reducing the density pa-

rameters, this would result in merging too many clusters together, making denser clusters indistinguishable.

Fuzzy Boundary Clusters. Some DCPGS geo-social clusters have fuzzy boundaries with each other, which is reasonable in the real world, since groups of socially connected users may spatially overlap. On the other hand, DBSCAN produces clusters with strict boundaries. For instance, in Figure 3(a), there is no strict boundary between the two clusters enclosed in region C. Although PureSocialDistance, which is the other extreme method, also produces clusters with fuzzy boundaries (see Figure 3(c)), the clusters are spatially indistinguishable and they are not interesting, i.e., for the applications mentioned in the Introduction.

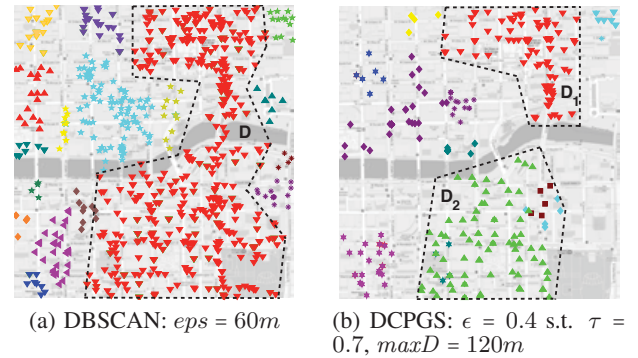


Figure 4: Clusters of Brightkite found by DBSCAN and DCPGS in Chicago

We also visually analyzed the results of the alternative graph-based clustering models LinkClustering and Metis, described in Section 2.1.1. These two approaches produce similar results; inductively, we show the clusters produced by LinkClustering in Figure 3(d). LinkClustering produces thousands of small clusters (average size around 3), which are typically not well-separated spa-

tially. Due to the sparsity of geo-social network data, the constructed place network contains many connected components that are disconnected with each other (e.g., the place network built when $\tau = 0.7$, $\max D = 100$, and $\omega = 0.5$ contains 34,496 connected components with 4.3 nodes and 8.2 edges on average. The clusters found by Metis are fewer and larger, but also spatially indistinguishable. Metis ignores outliers; as a result, places belong to the same cluster may have low spatial proximity and social similarity.

Finally, we analyzed the results of DCPGS, if our D_S definition (Definition 2) is replaced by the alternatives described in Section 2.2.1. For D_S^{Katz} and D_S^{ct} , we set $L = 3$; for bigger L values, these measures become extremely expensive. We observed that D_S^{Jac} , D_S^{Katz} , and D_S^{ct} produce similar results to each other. Indicatively, Figure 3(e) shows the clusters found by DCPGS if D_S^{Jac} is used instead of D_S . All these measures produce small clusters and too many outliers since they give large distance values for most pairs of places p_i and p_j . The set of common users for two places in Jaccard (i.e., $U_{p_i} \cap U_{p_j}$) is expected to be small and the decay factor of Katz dampens the effect of long connections between U_{p_i} and U_{p_j} . The expected CommuteTime distance between places is also high due to the effect of normalization. On the other hand, D_S^{sim} produces clusters of slightly larger sizes compared to D_S . We observed that the probability distribution of D_S^{sim} is skewed towards low values, meaning that many pairs of places have low bipartite minimax SimRank social distance, because SimRank is based on the most similar pair of visiting users. The clustering results of SimRank (Figure 3(f)) and DCPGS are visually similar; it is hard to tell which results are better based on visualization.

4.2 Social Quality Evaluation

In this section, we design and use two measures for assessing the social coherence between places in the discovered clusters. Based on these measures, we assess the quality of DCPGS and the alternative approaches for clustering GeoSN places.

4.2.1 Social Entropy based Evaluation

The first measure, called social entropy, measures the social quality of the clusters based on the network communities that the GeoSN users form. Given a social network $G = (U, E)$, we first partition all the users in U into several disjoint network communities. Let PC be a cluster of GeoSN places. According to the detected network communities, the visiting users of PC , i.e., $\cup_{p \in PC} U_p$, can be divided into several disjoint sets in $C_{PC} = \{C_1, C_2, \dots, C_m\}$, such that each set C_i is a subset of users in C_{PC} belonging to the same network community. We call C_{PC} the community set of PC .

Definition 3. (Social Entropy) Given a cluster PC , let U_{PC} be the set of users who visit the places in PC , i.e., $U_{PC} = \cup_{p \in PC} U_p$. The social entropy of PC is then defined as:

$$\mathcal{E} = \sum_{C_i \in C_{PC}} -\frac{|C_i|}{|U_{PC}|} \log \frac{|C_i|}{|U_{PC}|}$$

The social entropy, analogous to the entropy used in decision tree induction, measures the impurity of a cluster PC with respect to the participation of its users into different communities. A low social entropy means that the great majority of the visitors of PC come from the same community (i.e., low impurity), indicating that the places in PC have tighter social relationships between each other, which is favored.

We applied the METIS community detection algorithm [16] to divide the set of users in the GeoSN into k non-overlapping com-

munities⁷, providing a baseline for social entropy evaluation. To avoid a comparison that is biased to parameter k , we evaluate the social entropy of clustering results obtained by DCPGS and the competitors for two different values of k . One value of k is chosen based on the following *rule of the thumb* [19], i.e., $k = \sqrt{N}/2$ where N is the number of users in the social network. The other value of k is decided by Dunbar's number that suggests humans can only comfortably maintain 150 stable relationships and the mean community size is around 150. For dataset Gowalla, these values are $k = 313$ and $k = 1310$, respectively.

Figures 5(a) and 5(b) show the average social entropy for DCPGS, versions of DCPGS with alternative D_S (SimRank, CommuteTime, Katz, and Jaccard) and PureSocialDistance when varying ϵ . CommuteTime, and Katz have the lowest social entropy; however, these methods produce small clusters and have too many outliers as explained in Section 4.1. Within each small cluster, the places are only visited by few people and this explains the low entropy. Jaccard also has low social entropy for the same reason. PureSocialDistance has low social entropy in most cases; this indicates that our proposed social distance between places is effective in putting places with close social relationships together. When $\epsilon = 0.1$, the social entropies of DCPGS and PureSocialDistance are similar, both good, since only those places with very close social distances are clustered. When $\epsilon = 1$, PureSocialDistance has no social distance constraint τ thus its entropy becomes higher than that of DCPGS. DCPGS outperforms SimRank-based DCPGS, meaning that our proposed social distance is better than the SimRank-based $D_S^{sim}(p_i, p_j)$ distance. When $k = 1310$, the average social entropy of all the methods is larger than in the case where $k = 313$, since a larger number of network communities increases the probability that users in a single cluster belong to diverse communities. The average cluster size increases with ϵ , increasing the probability that the visitors of a cluster belong to different communities; thus, the average social entropy increases. In addition, the clustering result becomes stable at large values of ϵ , thus the social entropy converges. By visualization, we observed that ϵ should be set to a value smaller than 0.5 for the clustering results to be interesting.

Figures 5(c) and 5(d) show the average social entropies of DCPGS, DCPGS based on SimRank, CommuteTime, Katz, and Jaccard, and the graph-based clustering methods Metis and LinkClustering, when varying the social distance constraint τ . Similar to the case when ϵ varies, the social entropy increases and then stabilizes as τ increases, except for the entropy of Metis, which keeps increasing due to the network partitioning methodology of Metis with the increase of τ , the constructed place network becomes less connected, however, due to its partitioning nature, Metis puts disconnected places in the constructed place network into same cluster. When τ is less than 0.5, the social entropy of CommuteTime is zero, since with these distance measures the places in each cluster are visited by only one person when $\tau < 0.5$. Jaccard has low social entropy also due to the small sizes of its clusters. For $\tau \leq 0.1$ SimRank-based DCPGS fails to find any clusters, therefore the entropy is 0. After investigation, we found that there is no pair of places p_i, p_j with $D_S^{sim}(p_i, p_j) < 0.2$ because of the decay factor ϕ . When $\tau = 0.2$, SimRank has a low social entropy, since only few (987) clusters of small size are found compared to the 3605 clusters discovered by DCPGS. After the point where the two approaches find a similar number of clusters (e.g., at $\tau = 0.5$, SimRank finds 5880 clusters, while DCPGS finds 6742 clusters), DCPGS has constantly lower entropy than SimRank. In addition,

⁷This is different from Metis used as a competitor of DCPGS in GeoSN place clustering (discussed in Section 2.1.1).

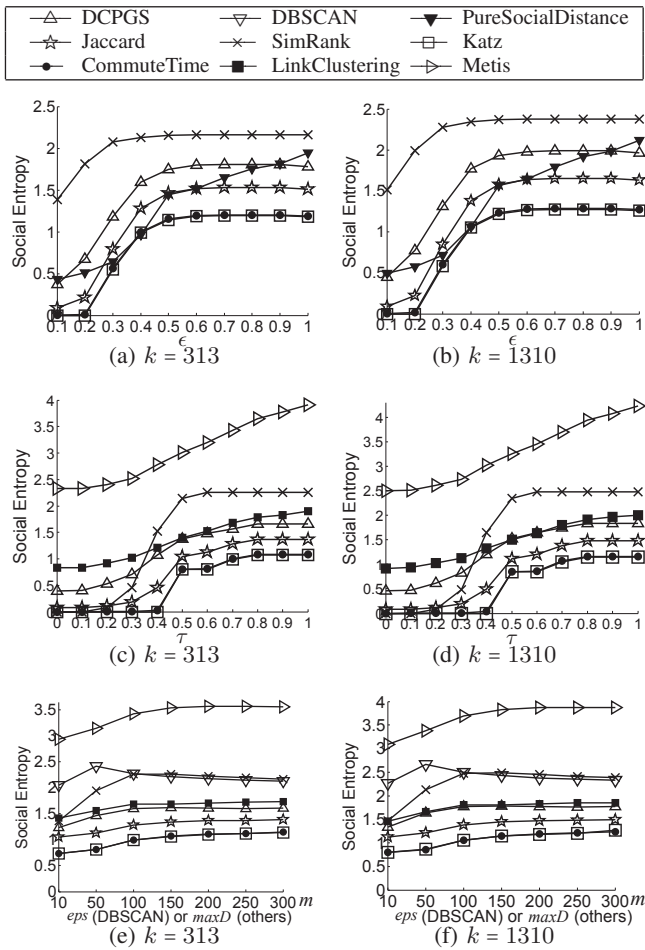


Figure 5: Social entropy evaluation in Gowalla

DCPGS is less sensitive to τ compared to SimRank. DCPGS outperforms the two graph-based competitors Metis and LinkClustering. As we observed by visualization, in practice τ should be set to a value higher than 0.5, because a very tight social distance constraint creates too few and too small geo-social clusters.

Figures 5(e) and 5(f) show the average social entropies of the various versions of DCPGS and all the competitor approaches when varying the spatial distance constraint $maxD$ (eps for DBSCAN). DCPGS is superior to SimRank-based DCPGS, DBSCAN, LinkClustering and Metis for all values of $maxD$ (eps). In general, the social entropies of all methods are not very sensitive to $maxD$. For Gowalla, a good value for $maxD$ is around 100m; large $maxD$ values result in clusters that are spatially too loose.

4.2.2 Community Score based Evaluation

Given a GeoSN place cluster PC , let U_{PC} be the set of users who visit the places in PC , i.e., $U_{PC} = \cup_{p \in PC} U_p$. Assume each U_{PC} is a community in the GeoSN. We adopt the eight network community multi-criterion scores surveyed in [18] to compute the community score of U_{PC} for each cluster PC . Figure 6 compares the results of DCPGS and its alternatives on Gowalla (Katz is omitted because its result is quite similar to that of CommuteTime), in terms of the *internal density* and *conductance* scores. We group the clusters discovered by each method by size and compute and plot the average community score (i.e., internal density and conductance) for each cluster size group. The results based on the other six criteria of [18] are similar and we omit them

due to lack of space. The internal density of U_{PC} is defined by $1 - m_{U_{PC}} / (|U_{PC}|(|U_{PC}| - 1)/2)$, where $m_{U_{PC}}$ is the number of edges in U_{PC} , $m_{U_{PC}} = |\{(u, v) | u \in U_{PC}, v \in U_{PC}\}|$. Conductance is the fraction of edges from nodes of U_{PC} that point outside U_{PC} , i.e., $o_{U_{PC}} / (2m_{U_{PC}} + o_{U_{PC}})$, where $o_{U_{PC}} = |\{(u, v) | u \in U_{PC}, v \notin U_{PC}\}|$. Let $f(U_{PC})$ be the community score of PC , based on either internal density or conductance; a smaller value of $f(U_{PC})$ indicates better social quality.

As Figures 6(a) and 6(c) show, the internal density increases with the cluster size. As the size of a place cluster increases, the denominator of the internal density formula increases quadratically while the number of social links between users in the cluster (i.e., the numerator) does not increase at the same pace. On the other hand, Figures 6(b) and 6(d) show that conductance initially decreases as the size of a cluster increases and fluctuates randomly for larger U_{PC} sizes, which is in line with the observations in [18]. In Figure 6, we observe that the geo-social clusters discovered by DCPGS have better community scores (i.e., lower internal density and conductance scores) compared to all competitors, except PureSocialDistance. Since PureSocialDistance uses our social distance in clustering and disregards spatial proximity, its social quality is expected to be better than that of DCPGS; still, as shown in Figure 3(c), its clusters are not distinguishable spatially. DCPGS outperforms DBSCAN and SimRank-based DCPGS. The quality gap between DCPGS and the 3 competitors in Figure 6(a) and 6(b) narrows as the size of clusters increases, since it is more difficult for a larger U_{PC} (usually obtained from a larger PC) to maintain a community-like structure compared to a smaller U_{PC} [18]. This indicates that our social distance is effective in finding geo-social clusters with small or medium size. DCPGS is also generally better than the four competitors in Figures 6(c) and 6(d). CommuteTime has better community scores when the cluster size is around 50. Most community scores of Jaccard, CommuteTime, LinkClustering, and Metis concentrate at the top-left corner of Figures 6(c) and 6(d), which indicates that these competitors have limited ability to discover geo-social clusters of various sizes. Furthermore, the quality gap between DCPGS and the five competitors in Figures 6(c) and 6(d) grows when the cluster size increases. We conclude that DCPGS (paired with our social distance measure) is the most effective method in finding geo-social clusters with both good social quality and identifiable spatial contour.

5. EFFICIENCY EVALUATION

In this section, we investigate the efficiency and scalability of DCPGS and its competitors. In Section 3, we presented two alternative implementations of DCPGS, an R-tree based (DCPGS-R) and a grid-based (DCPGS-G). We also implemented the corresponding R-tree based and grid-based versions of DBSCAN, and DCPGS based on Jaccard, SimRank, Katz, and CommuteTime. Linkclustering and Metis are also compared with DCPGS. In our comparison, we do not include PureSocialDistance, because this method is too slow (it takes roughly 1 day to finish) and it is not appropriate for spatial clustering.

Effect of data size. In the first experiment, we test the scalability of all methods. For this purpose, we use subsets of the two GeoSN datasets, Gowalla and Brightkite, of various sizes. The subsets are determined by continuous spatial ranges, which contain the desired cardinality of places. All place subsets are associated with the original social network data. Five subsets with 200K, 400K, 600K, 800K, and 1000K places are obtained from Gowalla. Similarly, five subsets of sizes 200K, 300K, 400K, 500K, and 600K are extracted from Brightkite. When varying the data size, we keep $\epsilon = 0.4$, $maxD = 100m$, $\tau = 0.7$ and $\omega = 0.5$ for all methods in

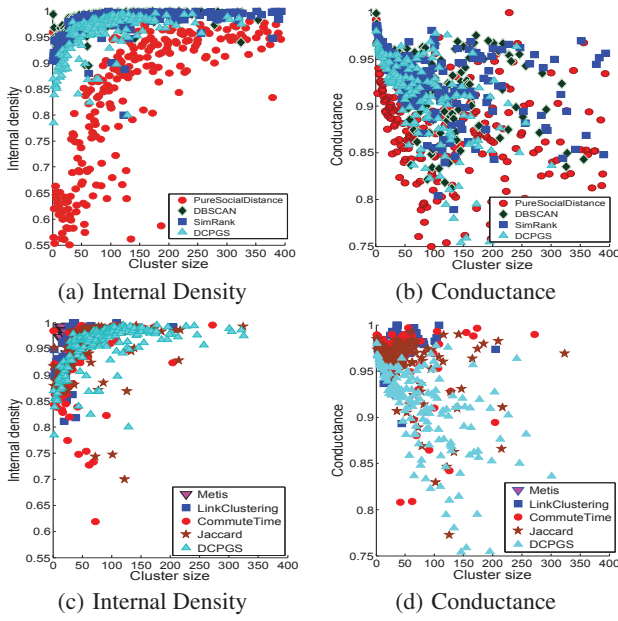


Figure 6: Community score evaluation in Gowalla

Gowalla and $\epsilon = 0.4$, $maxD = 120m$, $\tau = 0.7$ and $\omega = 0.5$ for all methods in Brightkite. (LinkClustering and Metis have no ϵ parameter.) Figures 7(a) and 7(b) display the runtime of all tested methods. DCPGS-G is much faster than DCPGS-R and maintains its performance gap as the database size grows. In all cases, DCPGS-G (DCPGS-R) is only a bit slower than DBSCAN-G (DBSCAN-R) and Jaccard-G (Jaccard-R)⁸, while DCPGS has much higher effectiveness as shown in Section 4. SimRank-G, (SimRank-R), Katz-G, and CommuteTime-G are much slower than DBSCAN-G (DBSCAN-R). This indicates that the time overhead imposed by DCPGS over the pure-spatial DBSCAN clustering algorithm is low, proving that our effective social distance is also efficient to compute. The costs of SimRank, Katz, and CommuteTime are all dominated by the expensive $D_S^{sim}(p_i, p_j)$, $D_S^{Katz}(p_i, p_j)$, and $D_S^{ct}(p_i, p_j)$ computations; optimizing their spatial search component (i.e., replacing the R-tree search by a grid-based search) has little effect on their overall performance.⁹ Metis has similar cost to DCPGS-G, while LinkClustering is slower; these methods are not effective for geo-social place clustering as discussed in Section 4.

Effect of ϵ . Figures 7(c) and 7(d) show the performance of all methods when varying ϵ , except LinkClustering and Metis, which do not use parameter ϵ . We keep $\tau = 0.7$ and $\omega = 0.5$ for both datasets and $maxD = 100m$ for Gowalla and $maxD = 120m$ for Brightkite. Note that DCPGS-G is faster than DCPGS-R in all cases, which again indicates that grid partitioning greatly improves the efficiency of DCPGS. DCPGS-G (DCPGS-R) is slightly more expensive than Jaccard-G (Jaccard-R). The runtimes of SimRank-G, SimRank-R, Katz-G, and CommuteTime-G increase significantly when ϵ increases from 0.1 to 0.5, because the increase weakens the power of the spatial filter (Proposition 1) and more (expensive) social distances ($D_S^{sim}(p_i, p_j)$, $D_S^{Katz}(p_i, p_j)$, and $D_S^{ct}(p_i, p_j)$) are computed. DCPGS-G and DCPGS-R remain

⁸To make the plots clearer, we omit the runtime of Jaccard-G (Jaccard-R) from some of them; the time of Jaccard-G (Jaccard-R) always falls between the time of DCPGS-G (DCPGS-R) and DBSCAN-G (DBSCAN-R).

⁹This is the reason why we omit the runtimes of Katz-R and CommuteTime-R from the comparison.

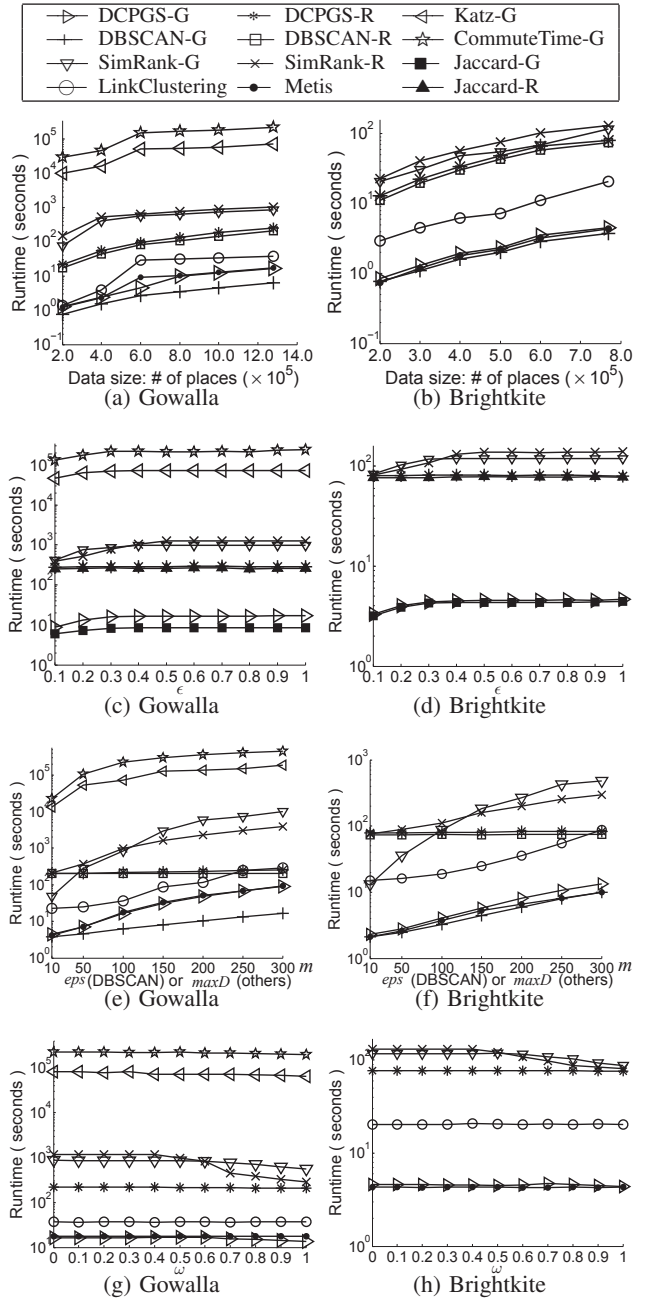


Figure 7: Runtime experiments

quite stable for all ϵ values since the computational cost of our social distance is insignificant.

Effect of $maxD$ (eps in DBSCAN). Parameter $maxD$ (eps in DBSCAN) decides the spatial ranges of the geo-social (spatial) neighborhoods of places and the size of the constructed place network for LinkClustering and Metis. Thus, $maxD$ (eps) greatly influences the efficiency of all methods. Figures 7(e) and 7(f) display the runtime of all methods when varying $maxD$ (eps). We set $\tau = 0.7$, $\epsilon = 0.4$, and $\omega = 0.5$ for all the methods except DBSCAN on both Gowalla and Brightkite datasets. Note that SimRank-G, SimRank-R, Katz-G and CommuteTime-G are more sensitive to $maxD$ compared to the other methods, due to the high cost of $D_S^{sim}(p_i, p_j)$, $D_S^{Katz}(p_i, p_j)$, and $D_S^{ct}(p_i, p_j)$ distance computations. The cost

of DCPGS-G stays below 100 seconds as $maxD$ ranges from 10m to 300m and remains close to that of DBSCAN-G and Jaccard-G. DCPGS-R is also close to DBSCAN-R and Jaccard-R. DCPGS-G remains faster than DCPGS-R, but their performance gap narrows with $maxD$; the reason is that grid-based computations become more expensive as the grid cell size (i.e., $maxD$) increases. DCPGS-R eventually becomes faster than DCPGS-G for a very large $maxD$, e.g. $maxD > 1000m$; however, as shown in Figure 3(c), geo-social clusters with a very large $maxD$ are not spatially identifiable. At an appropriate $maxD$ value (around 100m), DCPGS-G is clearly the best choice. Metis has similar cost to DCPGS-G while LinkClustering is slower.

Effect of ω . Figures 7(g) and 7(h) display the runtimes of the clustering methods when varying ω . We keep $\tau = 0.7$ and $\epsilon = 0.4$ for all methods on both datasets, while setting $maxD = 100m$ in Gowalla and $maxD = 120m$ in Brightkite. The costs of all methods are not much sensitive to ω ; an exception is the runtimes of SimRank-G, SimRank-R, Katz-G, and CommuteTime-G which decrease when $\omega > 0.5$. The reason is again the high effect of the spatial filter (Proposition 1), when ω is large, due to which many expensive social distance computations are avoided. As a final note, the social distance constraint τ is only applied as an ultimate filter and has no influence on the runtimes of all methods. It only slightly influences the costs of LinkClustering and Metis because it affects the size of the place network, on which these methods operate.

6. RELATED WORK

Our clustering problem is related to various research topics, including traditional spatial clustering, using mobility data to analyze places, clustering using spatial and non-spatial attributes, studying the relationship between spatial and social attributes, community detection, and other work on GeoSNs.

Spatial Clustering. Spatial clustering algorithms, surveyed in [12], are divided into three categories: *partitioning*, *hierarchical* and *density-based* clustering. Partitioning methods, including k -means, k -medoids, and CLARANS [23], are good at finding spherical-shaped clusters in small and medium-sized datasets. They need a pre-defined parameter k to specify the number of clusters obtained. However, partitioning methods are not able to detect clusters of arbitrary shapes. Hierarchical clustering techniques, such as BIRCH [41], Chameleon [15] and CURE [11], assign objects to clusters in two fashions: *agglomerative* (bottom-up) and *divisive* (top-down). Hierarchical clustering methods do not have well-defined termination criteria and cannot correct the result if some objects are assigned to the wrong clusters at an early stage. Density-based clustering methods, like DBSCAN [9], discover clusters of arbitrary shapes and sizes. Objects in dense regions are grouped as clusters, while objects in sparse regions are labeled as outliers. OPTICS [3] is an extension of DBSCAN, which generates an augmented ordering of the dataset that captures its density-based clustering structure at different granularities. DENCLUE [13] models the overall point density analytically as the sum of influence functions of the data points. Clusters can then be identified by determining density-attractors and clusters of arbitrary shapes can be easily described by a simple equation based on the overall density function. GDBSCAN [25] is a generalization of DBSCAN that clusters point objects as well as spatially extended objects according to both their spatial and their non-spatial attributes. Our work adopts the density-based clustering framework to find place clusters in a geo-social network, by considering both the spatial distance and the social coherence of the places.

Analysis of Places based on Mobility Data. Brillhante et al. [6] detect “communities” of places of interest (POIs) on a map based on how strongly the places are correlated in sequences of visits by mobile users. Different from our work, the social relationships between the users who visit the places and the spatial distances between places are disregarded. The co-existence of places in the visiting histories of users is the only criterion used for clustering. The proposed solution generates a graph G that connects pairs POIs according to the nature of their co-existence in sequences of user visits and then employs a classic algorithm for community detection on G to identify the place communities. Andrienko et al. [2] present an analysis and visualization tool, which first identifies interesting events of moving objects (e.g., instances of slow car movements), then spatially clusters these events, using density-based clustering to derive a set of significant places (e.g., regions where traffic jams occur), and finally applies visual analytics to aggregate and analyze the events with respect to parameters such as location, time and direction of movement.

Clustering based on Spatial and Non-Spatial Attributes. Clustering objects based on spatial and non-spatial attributes finds applications in different areas, such as computer vision, GIS, and social networks. Yu et al. [39] cluster pixels considering both the RGB color vectors and spatial proximity that is useful in natural image segmentation. Gennip et al. [32] use spectral clustering to identify communities in a graph where nodes are gang members and weighted edges indicate the gang members’ social interactions and geographic locations. Zhang et al. [40] apply clustering by adjusting the spatial distance between two objects according to the non-spatial attribute values between them.

Spatial-Social Relationship. The relationship between geography and social structure has been long studied by sociologists. Researchers have found that the likelihood of friendship with a person is decreasing with distance, which has been observed within colleges [29], new housing developments [10], and projects for the elderly [21]. Scellato et al. [27] performed a quantitative study on the socio-spatial properties of users in GeoSNs. By utilizing social and spatial properties of GeoSNs, the same research group proposed a link prediction model [28]. Ye et al. [37] designed a location recommendation system in the context of GeoSNs. Backstrom et al. [5] predict the location of an individual from a sparse set of known user locations using the relationship between geography and friendship. Wang et al. [34] find that the similarity between the movements of two individuals strongly correlates with their proximity in the social network. This correlation is used as a tool for link prediction in a social network. Pham et al. [24] propose an entropy-based model (EBM) that not only infers social connections but also estimates the strength of social connections by analyzing people’s co-occurrences in space and time. Different from existing work, we neither study the spatial-social relationship nor do prediction or recommendation utilizing this relationship. We perform density-based clustering of GeoSN places considering both the spatial distances between them and the social relationships between users that visit the places.

Detecting and Evaluating Communities in Networks. There are many existing works on network community detection and clustering of nodes in a graph using only the network distance between nodes [8, 22, 31, 35, 38]. SCAN [35] is an algorithm that detects clusters, hubs, and outliers in networks. [38] proposed partitioning, hierarchical and density-based algorithms to cluster objects on spatial networks, based on shortest-path distance. [18] summarized and empirically evaluated algorithms for network community detection. In Section 4.2.2, we have used network community quality

measures [18] to evaluate the social quality of the place clusters found by our algorithms.

Orthogonal Works in GeoSNSs. Cho et al. [7] studied the problem of using social relationships to explain human movement. Yang et al. [36] introduced a *Socio-Spatial Group Query* (SSGQ) that retrieves a group of users who have certain social connection strength and the sum of their spatial distances is minimized. Most recently, a general framework for Geo-Social query processing that supports *Range Friends* (RF), *Nearest Friends* (NF), and *Nearest Star Group* (NSG) queries was proposed in [4].

7. CONCLUSION

In this paper, we studied for the first time the problem of Density-based Clustering Places in Geo-Social Networks (DCPGS). Our clustering model extends the density-based clustering paradigm to consider both the spatial and social distances between places. We defined a new measure for the social distance between places, considering the social ties between users that visit them. Our measure is shown to be more effective and efficient to compute, compared more complex ones based on node-to-node graph proximity. We analyzed the effectiveness of DCPGS via case studies and demonstrated that DCPGS can discover clusters with interesting properties (i.e., barrier-based splitting, spatially loose clusters, clusters with fuzzy boundaries), which cannot be found by merely using spatial clustering. Besides, we designed two evaluation measures to quantitatively evaluate the social quality of clusters detected by DCPGS or competitors, called social entropy and community score, which also confirm that DCPGS is more effective than alternative approaches. Finally, we optimized the efficiency of DCPGS clustering with the help of a grid-based data structure. Our optimized DCPGS algorithm can cluster millions of places within just a few seconds. In the future, we plan to consider the influence of the time-span of user visits to places and the multiplicity of user visits to a place to cluster places.

8. REFERENCES

- [1] Y.-Y. Ahn, J. P. Bagrow, and S. Lehmann. Link communities reveal multiscale complexity in networks. *Nature*, 466:761–764, 2010.
- [2] G. L. Andrienko, N. V. Andrienko, C. Hurter, S. Rinzivillo, and S. Wrobel. Scalable analysis of movement data for extracting and exploring significant places. *IEEE TVCG*, 19(7):1078–1094, 2013.
- [3] M. Ankerst, M. M. Breunig, H.-P. Kriegel, and J. Sander. Optics: Ordering points to identify the clustering structure. In *SIGMOD*, 1999.
- [4] N. Armenatzoglou, S. Papadopoulos, and D. Papadias. A general framework for geo-social query processing. *PVLDB*, 6(10):913–924, 2013.
- [5] L. Backstrom, E. Sun, and C. Marlow. Find me if you can: Improving geographical prediction with social and spatial proximity. In *WWW*, 2010.
- [6] I. R. Brilhante, M. Berlingerio, R. Trasarti, C. Renso, J. A. F. de Macêdo, and M. A. Casanova. Cometogther: Discovering communities of places in mobility data. In *MDM*, 2012.
- [7] E. Cho, S. A. Myers, and J. Leskovec. Friendship and mobility: user movement in location-based social networks. In *KDD*, 2011.
- [8] A. Clauset, M. E. Newman, and C. Moore. Finding community structure in very large networks. *Physical review E*, 70(6):066111, 2004.
- [9] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, 1996.
- [10] L. Festinger, S. Schachter, and K. Back. Social pressures in informal groups: a study of human factors in housing. *Stanford Univ. Pr.*, 1963.
- [11] S. Guha, R. Rastogi, and K. Shim. Cure: An efficient clustering algorithm for large databases. In *SIGMOD*, 1998.
- [12] J. Han and M. Kamber. *Data Mining: Concepts and Techniques*. Morgan Kaufmann, 2000.
- [13] A. Hinneburg and D. A. Keim. An efficient approach to clustering in large multimedia databases with noise. In *KDD*, 1998.
- [14] G. Jeh and J. Widom. Simrank: A measure of structural-context similarity. Technical Report 2001-41, Stanford InfoLab, 2001.
- [15] G. Karypis, E.-H. Han, and V. Kumar. Chameleon: Hierarchical clustering using dynamic modeling. *IEEE Computer*, 32(8):68–75, 1999.
- [16] G. Karypis and V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on scientific Computing*, 20(1):359–392, 1998.
- [17] J. Leskovec, J. M. Kleinberg, and C. Faloutsos. Graph evolution: Densification and shrinking diameters. *TKDD*, 1(1), 2007.
- [18] J. Leskovec, K. J. Lang, and M. W. Mahoney. Empirical comparison of algorithms for network community detection. In *WWW*, 2010.
- [19] K. V. Mardia, J. Kent, and J. Bibby. *Multivariate analysis (probability and mathematical statistics)*. Academic Press London, 1980.
- [20] S. Milgram. The small world problem. *Psychology today*, 2(1):60–67, 1967.
- [21] L. Nahemow and M. Lawton. Similarity and propinquity in friendship formation. *Journal of Personality and Social Psychology*, 32(2):205–213, 1975.
- [22] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, 2001.
- [23] R. T. Ng and J. Han. Efficient and effective clustering methods for spatial data mining. In *VLDB*, 1994.
- [24] H. Pham, C. Shahabi, and Y. Liu. Ebm: an entropy-based model to infer social strength from spatiotemporal data. In *SIGMOD*, 2013.
- [25] J. Sander, M. Ester, H.-P. Kriegel, and X. Xu. Density-based clustering in spatial databases: The algorithm gbscan and its applications. *Data Min. Knowl. Discov.*, 2(2):169–194, 1998.
- [26] P. Sarkar and A. W. Moore. A tractable approach to finding closest truncated-commute-time neighbors in large graphs. In *UAI*, 2007.
- [27] S. Scellato, A. Noulas, R. Lambiotte, and C. Mascolo. Socio-spatial properties of online location-based social networks. In *ICWSM*, 2011.
- [28] S. Scellato, A. Noulas, and C. Mascolo. Exploiting place features in link prediction on location-based social networks. In *KDD*, 2011.
- [29] J. Q. Stewart. An inverse distance variation for certain social influence. *Science*, 93(2404):89–90, 1941.
- [30] P.-N. Tan, M. Steinbach, and V. Kumar. *Introduction to Data Mining*. Addison-Wesley, 2005.
- [31] C. Tantipathananandh, T. Y. Berger-Wolf, and D. Kempe. A framework for community identification in dynamic social networks. In *KDD*, 2007.
- [32] Y. van Gennip, B. Hunter, R. Ahn, P. Elliott, K. Luh, M. Halvorson, S. Reid, M. Valasik, J. Wo, G. E. Tita, A. L. Bertozzi, and P. J. Brantingham. Community detection using spectral clustering on sparse geosocial data. *CoRR*, abs/1206.4969, 2012.
- [33] C. Wang, V. Satuluri, and S. Parthasarathy. Local probabilistic models for link prediction. In *ICDM*, 2007.
- [34] D. Wang, D. Pedreschi, C. Song, F. Giannotti, and A.-L. Barabasi. Human mobility, social ties, and link prediction. In *KDD*, 2011.
- [35] X. Xu, N. Yuruk, Z. Feng, and T. A. J. Schweiger. Scan: a structural clustering algorithm for networks. In *KDD*, 2007.
- [36] D.-N. Yang, C.-Y. Shen, W.-C. Lee, and M.-S. Chen. On socio-spatial group query for location-based social networks. In *KDD*, 2012.
- [37] M. Ye, P. Yin, and W.-C. Lee. Location recommendation for location-based social networks. In *GIS*, 2010.
- [38] M. L. Yiu and N. Mamoulis. Clustering objects on a spatial network. In *SIGMOD*, 2004.
- [39] Z. Yu, O. C. Au, R. Zou, W. Yu, and J. Tian. An adaptive unsupervised approach toward pixel clustering and color image segmentation. *Pattern Recognition*, 43(5):1889–1906, 2010.
- [40] B. Zhang, W. J. Yin, M. Xie, and J. Dong. Geo-spatial clustering with non-spatial attributes and geographic non-overlapping constraint: A penalized spatial distance measure. In *PAKDD*, 2007.
- [41] T. Zhang, R. Ramakrishnan, and M. Livny. Birch: An efficient data clustering method for very large databases. In *SIGMOD*, 1996.