

距离-关键字相似度约束的双色反 k 近邻查询方法

张豪*, 朱睿, 宋楸尧, 方鹏, 夏秀峰

(沈阳航空航天大学 计算机学院, 沈阳 110136)

(* 通信作者电子邮箱 13623212292@163.com)

摘要:针对空间关键字双色反 k 近邻查询返回结果质量较低的问题,提出了基于距离-关键字相似度约束的双色反 k 近邻查询方法。首先,通过设置一个阈值将查询结果中质量较低的用户给过滤掉,从而避免了查询结果中出现空间距离相对较远的用户,保证了查询结果质量;然后,为支持该查询,提出了一种关键字多分辨率网格矩形树(KMG-Tree)索引来管理数据;最后,提出了基于Six-region算法的Six-region-optimize算法来提高查询处理效率。Six-region-optimize算法的查询效率相较baseline和Six-region算法分别平均提高了约85.71%和23.45%。基于真实时空数据进行实验测试和分析,实验结果验证了Six-region-optimize算法的有效性和高效性。

关键词:关键字;双色反 k 近邻查询;空间距离;相似度约束;查询效率

中图分类号: TP311 文献标志码: A

Bichromatic reverse k nearest neighbor query method based on distance-keyword similarity constraint

ZHANG Hao*, ZHU Rui, SONG Fuyao, FANG Peng, XIA Xiufeng

(College of Computer Science, Shenyang Aerospace University, Shenyang Liaoning 110136, China)

Abstract: In order to solve the problem of low quality of results returned by spatial keyword bichromatic reverse k nearest neighbor query, a bichromatic reverse k nearest neighbor query method based on distance-keyword similarity constraint was proposed. Firstly, a threshold was set to filter out the low-quality users in the query results, so that the existence of users with relatively long spatial distance in the query results was avoided and the quality of the query results was ensured. Then, in order to support this query, an index of Keyword Multiresolution Grid rectangle-tree (KMG-tree) was proposed to manage the data. Finally, the Six-region-optimize algorithm based on Six-region algorithm was proposed to improve the query processing efficiency. The query efficiency of the Six-region-optimize algorithm was about 85.71% and 23.45% on average higher than those of the baseline and Six-region algorithms respectively. Experimental test and analysis were carried out based on real spatio-temporal data. The experimental results verify the effectiveness and high efficiency of the Six-region-optimize algorithm.

Key words: keyword; bichromatic reverse k nearest neighbor query; spatial distance; similarity constraint; query efficiency

0 引言

随着移动互联网技术的不断发展,基于位置的社交网络(Location-Based Social Network, LBSN)服务应用越来越广泛。许多学者研究了面向LBSN的查询处理问题应对不同服务类型下的业务请求。

在众多查询处理问题中,面向反 k 近邻空间关键字(Reverse Spatial Keyword k Nearest Neighbor, RSK k NN)^[1]查询是一类重要问题,在市场分析、决策支持和交通信息等领域具有重要应用。具体地,在客户推荐系统中,该查询可以帮助商家根据客户偏好有针对性地利用短信等手段为客户推荐商品信息,从而增加营业收入。在交通信息管理领域,它可以通过分析乘客的用车偏好和位置信息,为乘客推荐车辆。

传统反 k 近邻查询可分为单色反 k 近邻查询(简称单色查询)和双色反 k 近邻查询(简称双色查询)。具体地,给定一组对象 o 和一个单色查询 q ,查询返回所有以 q 作为其 k 近邻的对象。与之不同,给定一组用户 U 、一组设施 F 和一个双色查询设施 q ,查询返回以设施 q 为 k 近邻的所有用户。反 k 近邻查询主要应用于推荐场景,其中单色反 k 近邻查询主要应用于同类型对象之间的推荐。例如,为某位用户推荐一些与他距离比较近的志趣相投的人。而双色反 k 近邻主要应用于两种不同类型对象之间的推荐,相较于单色反 k 近邻查询的应用更加广泛,比如推荐场景中更多的是为商场提供潜在的消费者,为司机提供潜在的乘客等,这些都是单色反 k 近邻查询无法做到的,所以本文重点研究双色反 k 近邻。接下来本文

收稿日期: 2020-09-18; 修回日期: 2020-12-09; 录用日期: 2020-12-10。 基金项目: 国家自然科学基金资助项目(61702344)。

作者简介: 张豪(1996—),男,河北石家庄人,硕士研究生,主要研究方向:大数据; 朱睿(1982—),男,辽宁沈阳人,副教授,博士,CCF会员,主要研究方向:大数据; 宋楸尧(1995—),男,辽宁鞍山人,硕士研究生,主要研究方向:大数据; 方鹏(1994—),男,甘肃平凉人,硕士研究生,主要研究方向:新媒体大数据; 夏秀峰(1964—),男,山东胶南人,教授,博士,主要研究方向:数据库、数据仓库、数据挖掘。

以 $k=2$ 为例简单说明传统反 k 近邻的查询规则,如图1(a)所示,对象 $\{o_0, o_1, o_5, o_6\}$ 的 k 近邻对象包含 q ,所以对对象 q 的单体反 k 近邻查询返回结果为 $\{o_0, o_1, o_5, o_6\}$;如图1(b)所示,用户 $\{u_0, u_5, u_6, u_7, u_9\}$ 的 k 近邻设施包含 q ,所以设施 q 的双色反 k 近邻查询结果为 $\{u_0, u_5, u_6, u_7, u_9\}$ 。本文研究双色反 k 近邻查询。

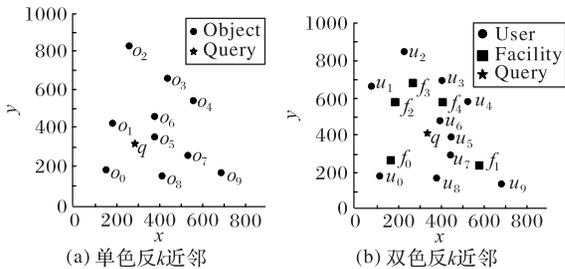


图1 传统单色和双色反 k 近邻

Fig. 1 Traditional monochromatic and bicolor reverse k nearest neighbors

传统反 k 近查询不能根据设施和用户需求针对性地返回查询结果,而带关键字约束的反 k 近邻查询可以很好解决此问题。鉴于此类问题的重要性,许多学者针对此类问题展开研究。例如:Lu等^[1]首次研究了基于关键字约束的反 k 近邻查询问题。给定查询点 q 和对象集合 D ,查询计算 q 和 D 中对象空间上的距离以及关键字间的相似性。以此为基础,算法利用打分函数评价各对象对于查询点 q 的重要性,并返回得分最高的 k 个对象。Zhao等^[2]提出了一种相关反近邻布尔空间关键字查询(Ranked Reverse boolean Spatial Keyword Nearest Neighbors, Ranked-RSKNN)问题。该查询返回一定数量与查询点相关程度最高的对象。

然而,这些研究存在的共性问题是没有考虑查询点与查询结果的距离对查询结果造成的影响,很多距离查询点很远但是关键字相似度较高的对象也将返回给用户。显然,和一些距离查询点较近但关键字相似度较低的对象相比,这些查询结果的质量较低。

因此,本文研究距离-关键字相似度约束的双色反 k 近邻(Distance-Keyword Reverse k Nearest Neighbor, DKR k NN)问题。和以往研究不同,DKR k NN通过引入阈值 φ ,过滤掉一些距离查询点相对较远的用户。为高效支持此类查询,本文需面对以下挑战:

1)高效的过滤能力。空间文本数据同时具有空间位置和关键字属性。因此,如何同时根据空间位置和关键字对数据进行筛选是极具挑战的。

2)高效的验证能力。对每个用户进行验证都需要遍历一次设施集。因此,在验证时如何尽可能少地遍历设施是极具挑战的。

针对上述挑战,本文提出了一种基于多分辨率网格树查询处理框架支持查询。本文主要工作如下:

1)提出了一种关键字多分辨率网格矩形树(Keyword Multiresolution Grid rectangle-Tree, KMG-Tree)索引管理设施和用户数据。该索引同时具有以下优点:①通过对网格中的单元格添加关键字信息,以达到通过这两个属性批量过滤数

据的目的;②相对平衡的结构,通过对传统的多分辨率网格进行改进,KMG-Tree用分辨率较高网格划分数据密集的区域;用分辨率较低网格划分数据稀疏的区域。这样一来,它可以保证索引各节点维护的对象数目大致相同从而保证了索引的平衡性。

2)提出了基于Six-region-optimize的过滤算法。该算法将整个空间等分为6个区域,根据位置关系定义过滤区域。这样一来,被过滤区域覆盖的网格节点可被安全过滤。此外,与查询点关键字无交集的网格节点也可被安全过滤。

3)提出了基于Six-region-optimize的验证算法。在验证用户时,只需要访问该区域及相邻区域的对用户可能有影响的设施,这样可以快速找到用户的最近邻及 k 近邻,一定程度上提高了验证的效率。

1 相关工作

LBSN平台的数据规模越来越庞大,在做大数据查询时首先要对大数据进行分析,例如对于大数据分类问题,Xia等^[3]提出了粒度球邻域粗糙集(Granular Ball Neighborhood Rough Set, GBNRS)分类方法,该方法自适应地为每个对象生成不同的邻域,从而具有更大的通用性和灵活性,并提高了在公共基准数据集上的分类精度和性能。对于大数据“ k -均值”问题,Xia等^[4]提出了一种称为“Ball k -means”的方法,用球来描述每个簇,减少了点质心距离的计算,提高了算法的性能。对于大数据噪声监测问题,Xia等^[5]提出了一种完整有效的随机森林方法,通过模拟网格的生成和扩展来检测类噪声,提高了噪声监测的性能。在这些大数据分析的基础上许多学者根据大数据推荐问题研究了反 k 近邻查询。

区域剪枝是处理空间反 k 近邻查询的方法之一,其主要思想是将空间分为几个区域,并针对每块区域的空间距离对用户进行剪枝。Stanoi等^[6]提出了一种基于Six-region剪枝的方法。该方法以查询点为中心将空间分为6个相等的区域。在每个区域中找到查询点的 k 近邻设施,以该设施到查询点距离为半径定义剪枝区域。

半空间剪枝法的核心思想是找到查询对 q 与设施对象 f 的中垂线。以此为基础,Tao等^[7-8]提出了半空间剪枝的算法,给定一个设施 f 和一个查询 q , f 和 q 之间的垂直平分线 $B_{f,q}$ 将空间分成两部分。让 $H_{f,q}$ 表示包含 f 的半空间, $H_{q,f}$ 表示包含 q 的半空间。 $H_{f,q}$ 中的每个用户 u 满足 $dist(u, f) < dist(u, q)$ 。换句话说,位于 $H_{f,q}$ 中的每个用户 u 被 f 剪枝,被 k 个这样的设施剪枝的用户可以被过滤。Cheema等^[9-11]提出影响区概念,影响区是指当且仅当 u 位于该区域内时, u 为 q 的R k NN。建造影响区的方法是绘制所有设施的半空间,被小于 k 个设施过滤的区域为影响区。Lee等^[12]提出了相关反最近邻查询,该查询通过对数据对象与查询对象 q 的相关性进行排序,返回受查询对象 q 影响最大的 t 个数据对象。

Yiu等^[13]提出了在大型图上的反 k 近邻查询,并提出了两种剪枝方法:一种是在网络节点被访问后立即进行修剪的急切算法;另一种是在发现数据点时对搜索空间进行修剪的懒惰算法,实现了在大型图上的反 k 近邻查询。李佳佳等^[14]提出了面向时间依赖的路网中的反 k 近邻查询,该查询提出了

计算节点到达时间的方法,利用该到达时间查询出多个候选 k 近邻结果。Hidayat 等^[15]提出了近似反最近邻查询,该查询将一些不是查询点的反最近邻但距离查询点也很近的点返回到了查询结果中。

以上几个算法都只是在空间或路网中反 k 近邻查询应用很好的算法,并不能够解决反 k 近邻空间关键字查询问题。Lu 等^[1]首次提出了反 k 近邻空间关键字查询问题,将空间位置和关键字通过度量公式结合在一起,并设计了一种空间文本混合索引结构 IUR-Tree (Intersection-Union R-Tree)。Zhao 等^[2]提出了一种相关反最近邻空间关键字查询,简称 Ranked-RS k NN。该查询可以返回固定数量的对象,并将这些对象按照查询相关程度的大小进行排序。该查询能保证找到一定数量的结果返回。

然而,在反 k 近邻空间关键字查询问题中,没有考虑空间距离太远和关键字相似度过低的情况,得到的查询结果并非是查询点所期望的结果。因此,研究基于距离-关键字相似度约束的反 k 近邻空间关键字查询问题具有重要意义。

2 问题定义

前文已经提到现有的空间关键字双色反 k 近邻查询都没有考虑查询点与查询结果之间的距离对查询结果造成的影响,很多距离查询点很远但是关键字相似度较高的对象也将返回给用户。显然,和一些距离查询点较近但关键字相似度较低的对象相比,这些查询结果的质量较低。具体地,在商场推销商品的场景下,通过双色反 k 近邻查询找到了一些潜在的消费者。当 $k \geq 2$ 时,这些消费者可能存在这样一种情况,商场 a 的位置在城市的中心,但是它的反 k 近邻中存在某个消费者 u 在该城市的郊区,在该郊区存在另一个商场 b 也可以满足 u 的消费需求,商场 a 对于 u 来说太远了,所以消费者 u 一定不会舍近求远,那么查询返回的消费者 u 对于商场 a 来说质量是非常低的。本文针对该问题设置了一个阈值 φ_r ,通过阈值 φ_r 可以将这些质量非常低的对象给过滤掉。下面详细介绍本文查询问题的形式化定义。

在介绍问题之前,本文首先定义一组符号。给定用户集合 U 和设施集合 F ,任意用户 $u \in U$ 、任意设施 $f \in F$ 均可用二元组 $\langle l, key \rangle$ 表示。其中, $u.l$ 表示了 u 的位置信息, $u.key$ 表示了用户的关键字信息。给定任意查询 q ,它可用五元组 $\langle l, key, k, \varphi_r, \varphi_k \rangle$ 表示。 φ_r 表示用户到查询设施距离与其到最近邻设施距离比值的阈值; φ_k 表示用户与查询设施关键字相似度阈值。给定用户 u 和设施 f ,本文利用欧氏距离 $dist(u, f)$ 计算 u 与 f 之间的距离。给定 u 和其最近邻 $NN(f)$,本文利用式(1)计算 φ_r ,利用式(2)计算 u 与 q 所对应关键字的相似程度。接下来介绍关键字-距离约束条件下的双色反 k 近邻查询。

$$Ratio(q, u) = \frac{dist(q, u)}{dist(u, NN(f))} \quad (1)$$

$$Sim(q, u) = \frac{(q.k \cap u.k)}{u.k} \quad (2)$$

定义 1 距离-关键字相似度约束的双色反 k 近邻查询 (DKR k NN)。给定一组带关键字的用户 U 、一组带关键字的设

施 F 、一个带关键字的查询设施 q 、用户到查询设施距离与其最近邻设施距离的比值阈值 φ_r 和关键字相似度阈值 φ_k ,返回满足条件 $Ratio(q, u) \leq \varphi_r$ 和 $Sim(q, u) \geq \varphi_k$,并以查询设施 q 为 k 近邻的每个用户 u 。

接下来,通过如下实例详细说明 DKR k NN 查询。通过设施和用户来区分不同的数据类型。图 1(b)为每个对象在二维空间中的位置。图中五角星 q 表示商场在空间中的位置,正方形 f_i 表示其他商场在空间中的位置,圆形 u_i 表示用户在空间中的位置(后文图中的正方形点代表设施,圆形点代表用户)。表 1 分别为用户以及用户与商场 q 的关键字相似度、用户到商场 q 的空间距离。给定查询 q ,假设查询值 $k = 2$ 、 $\varphi_r = 2$ 、 $\varphi_k = 0.3$,关键字集合为 $\{\text{water, egg, meat, coffee, banana, milk}\}$, q 的 R2NN 为 $\{u_0, u_5, u_6, u_7, u_9\}$ (见图 1(b))。 u_0 的最近邻为 f_0 , u_9 的最近邻为 f_1 ,通过计算发现 $Ratio(q, u_0) = 3.05 > 2$, $Ratio(q, u_9) = 3.47 > 2$ (见表 1)。因此, u_0 和 u_9 不满足要求。进一步地,通过相似度计算发现用户 $\{u_5, u_6, u_7\}$ 的关键字信息与商场 q 的关键字的信息相似度分别为 $\{0.09, 0.71, 0.50\}$,显然 $Sim(q, u_5) < 0.3$,最终,查询结果为 $\{u_6, u_7\}$ 。

表 1 用户与设施之间的关系

Tab. 1 Relationship between users and facilities

用户	相似度	到查询点的距离	最近邻设施距离
u_0	0.50	$dist(q, u_0) = 220.91$	$dist(f_0, u_0) = 72.3$
u_1	0.63	$dist(q, u_1) = 360.69$	$dist(f_2, u_1) = 142.13$
u_2	0.25	$dist(q, u_2) = 451.22$	$dist(f_3, u_2) = 167.63$
u_3	0.80	$dist(q, u_3) = 288.62$	$dist(f_4, u_3) = 100.05$
u_4	0.36	$dist(q, u_4) = 261.73$	$dist(f_4, u_4) = 110.00$
u_5	0.09	$dist(q, u_5) = 118.80$	$dist(q, u_5) = 118.80$
u_6	0.71	$dist(q, u_6) = 92.20$	$dist(q, u_6) = 92.20$
u_7	0.50	$dist(q, u_7) = 161.55$	$dist(f_1, u_7) = 148.66$
u_8	0.21	$dist(q, u_8) = 225.61$	$dist(f_1, u_8) = 206.16$
u_9	0.55	$dist(q, u_9) = 424.38$	$dist(f_1, u_9) = 122.07$

3 基于多分辨率网格树的查询处理框架

3.1 框架概述

本文利用基于 Six-region 改进的 Six-region-optimize 过滤-验证框架支持基于距离-关键字相似度约束的双色反 k 近邻查询问题。为支持高效查询,本文首先提出了一种基于数据依赖的多分辨率网格树索引实现数据的高效管理,随后本文提出了一种基于 Six-region-optimize 过滤算法实现对用户和设施的过滤,最后本文提出了一种基于 Six-region-optimize 的验证算法实现对候选用户的验证。

给定 n 个带关键字的用户,首先,根据用户空间位置将空间等分为 $\sqrt{n} \times \sqrt{n}$ 个带关键字信息的单元格。然后,根据单元格中数据的密度对单元格进行合并和分裂操作,从而得到一个多分辨率网格结构。最后,通过多叉树维护网格之间的空间位置关系。该索引可以利用网格的特性降低维护代价,保证索引的平衡性,并在访问过程中根据单元格的关键词信息对用户实现批量过滤。

给定查询 $q \langle q.l, q.k \rangle$,过滤算法以 q 为中心将整个空间划分为6个区域 $\{S_0, S_1, \dots, S_5\}$ 。对于每一个区域,算法利用用户和设施之间空间位置关系与关键字信息分别对设施和用户进

行过滤,得到候选用户集合 R 和用于验证用户的设施集合 V 。这样做的好处是可以快速过滤掉一些用户和设施,从而减少计算代价。

根据过滤算法得到候选用户集 R 和用于验证用户的设施集 V ,根据验证算法对候选集 R 中的每个用户进行验证。验证每个候选用户需要遍历设施集 V ,将满足查询要求的用户存储到结果集 $Result$ 中。

3.2 多分辨率网格树索引 KMG-Tree

本文提出了基于关键字多分辨率网格矩形树(KMG-Tree)索引分别管理用户和设施。该索引是一种以多分辨率网格为基础的数据结构。它对数据密集的区域使用分辨率较高的单元格进行划分,对数据稀疏的区域使用分辨率较低的单元格进行划分,从而保证每个单元格存储用户或设施数量基本相同。

首先,介绍一组 KMG-Tree 索引用到的符号。KMG-Tree 索引的每个节点 t 使用一个 5 元组 $\langle id, h, c, k, leaf \rangle$ 表示。KMG-Tree 每个节点 t 唯一标识用 $t.id$ 表示; $t.c$ 是一个数组,用来存储节点 t 的所有孩子节点; $t.leaf$ 是一个标记,它标记节点 t 是否为 KMG-Tree 索引的叶子节点,值为 true 表示为该索引的叶子节点,值为 false 表示为该索引的非叶子节点; $t.k$ 用来存储每个节点中所有对象的关键字信息,该关键字信息是一个基于哈夫曼树构造的一个 bit 串。算法可以根据每个节点 t 存储的 bit 串快速过滤对象。使用哈夫曼编码表示关键字信息可以有效地降低空间代价。与直接使用节点存储关键字方式相比,基于哈夫曼编码的 bit 串存储策略可以有效压缩存储空间。

接下来,举一个例子详细说明索引的构建过程。本文首先指定一个阈值 m 用来表示每个单元格存储对象数量的最大值。以 $m = 2$ 时为例,整个空间存在 10 个对象,将整个空间划分为 $\sqrt{10} \times \sqrt{10}$ 个单元格,根据每个单元格中对象的数量判断单元格是否需要合并或者分裂。合并和分裂的条件如下:

- 1)当单元格中对象数量大于 m 时进行分裂;
- 2)当单元格中对象数量小于 m 时进行合并。

经过合并与分裂之后得到的逻辑空间如图 2(a)所示, C_0 包含四个节点 C_1, C_2, C_3 和 C_4 ; C_1 包含两个对象 u_1 和 u_2 ; C_2 包含两个节点 C_5 和 C_6 ; C_3 包含一个对象 u_0 ; C_4 包含两个节点 C_7 和 C_8 ; C_5 包含一个对象 u_3 ; C_6 包含两个对象 u_4 和 u_6 ; C_7 包含两个对象 u_5 和 u_7 ; C_8 包含两个对象 u_8 和 u_9 。这时所有单元格都满足对象数量的限制。将 C_1, C_3, C_5, C_6, C_7 和 C_8 的 $leaf$ 值设置为 true。构建完成后的 KMG-Tree 索引结构如图 2(b)所示,索引中每个节点以及每个单元格集成对应的关键字集合存储在 $t.k$ 中。

接下来,介绍构建 KMG-Tree 时间复杂度。构建 KMG-tree 索引需要遍历所有的对象,构建用户索引首先遍历所有用户时间代价为 $O(|U|)$,遍历过程需要将用户存储到网格节点中。根据网格划分得到 $\sqrt{|U|} \times \sqrt{|U|}$ 个单元格,索引的高度为 $\log \sqrt{|U|} \times \sqrt{|U|}$,分裂与合并的时间代价为 $O(\log \sqrt{|U|} \times \sqrt{|U|})$,所以构建用户 KMG-Tree 索引的总代价为 $O(|U| +$

$\log \sqrt{|U|} \times \sqrt{|U|})$ 。同理构建设施 KMG-Tree 索引的总代价为 $O(|V| + \log \sqrt{|V|} \times \sqrt{|V|})$ 。

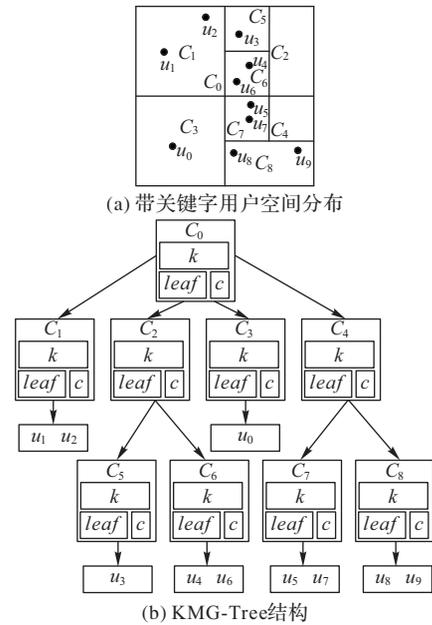


图 2 索引结构示意图

Fig. 2 Schematic diagram of index structure

3.3 基于 Six-region-optimize 的过滤算法

Six-region 算法在处理双色反 k 查询时,不仅没有考虑关键字因素,而且只是对用户进行了过滤,这在验证用户时需要遍历整个设施 R-tree。针对以上问题本文提出了基于 Six-region 改进的过滤算法 Six-region-optimize,该算法结合设施空间位置与关键字结合对用户和设施进行过滤。在介绍算法之前先介绍一个引理和一个定理。

引理 1 以查询点 q 为中心将空间划分为 6 个相等区域,在区域 $i(i = 0, 1, \dots, 5)$ 中,找到距离设施 q 第 k 近的设施 f_k 。以 q 为圆心、 $r_i = dist(q, f_k)$ 为半径做一个扇形 sec_i ,位于 sec_i 外面的用户可以被过滤。

证明 在区域 $i(i = 0, 1, \dots, 5)$ 中,位于 sec_i 外面的任意用户 u 的 kNN 一定不包含 q ,因为在 sec_i 里面一定存在 k 个设施满足 $dist(u, f) < dist(u, q)$ 。

定理 1 以查询点 q 为中心将空间划分为 6 个相等区域,在区域 $i(i = 0, 1, \dots, 5)$ 中,以 q 为圆心、 $2R_i$ 为半径做一个扇形 SEC_i ,位于 SEC_i 外面的设施可以被过滤。其中 R_i 为区域 i 及其相邻区域中 r 的最大值。 r 是根据引理 1 得到的每个区域的 sec_i 的半径。

证明 在区域 $i(i = 0, 1, \dots, 5)$ 中,位于 SEC_i 外面的任意设施 f 一定不会影响查询点 q 成为任何候选用户 u 的 kNN 。因为每个区域中的任意候选用户 u 一定满足 $dist(q, u) \leq R_i$,而位于 i 区域 SEC_i 外面的任意设施 f 一定满足 $dist(f, u) \geq R_i$ 。

接下来,介绍基于 Six-region-optimize 的过滤算法。该算法过滤过程分为批量过滤和局部过滤。针对用户的批量过滤,在遍历用户 KMG-Tree 过程中根据定理 1 对用户空间位置上进行过滤,在区域 $i(i = 0, 1, \dots, 5)$ 中位于 sec_i 外面的网格节点可以直接被过滤。然后,根据关键字信息对用户进行过

滤,没有被直接过滤的网格节点 t 需要检查该节点的关键字信息,如果 t 满足 $t.k \cap q.k = \emptyset$,则 t 可以被直接过滤。针对用户的局部过滤,访问那些没有被直接过滤网格节点的用户,根据定理1过滤掉那些在区域 $i(i=0,1,\dots,5)$ 中位于 sec_i 外面的用户。

例如图3(a)所示,以查询 q 为中心的整个空间被划分为六个相等的扇形区域 $S_0 \sim S_5$ 。在区域 S_1 和 S_2 中,查询点 q 的2近邻设施分别为 f_4 和 f_0 ,所以基于这两个设施定义了这两个区域的过滤空间,如图中阴影部分。被阴影部分完全覆盖的网格节点可以直接被过滤;被部分覆盖的网格节点将访问所包含的网格节点或用户,并将位于过滤空间中的网格节点或用户过滤。这样一来,算法无需要遍历所有用户便可批量过滤掉不满足要求的用户。

针对设施的批量过滤,在遍历设施KMG-Tree过程中根据定理1对设施在空间位置上进行过滤,在区域 $i(i=0,1,\dots,5)$ 中位于 SEC_i 外面的网格节点可以直接被过滤。针对设施的局部过滤,访问那些没有被直接过滤网格节点的设施,根据定理1过滤掉那些在区域 $i(i=0,1,\dots,5)$ 中位于 SEC_i 外面的设施。

例如图3(b)所示,在区域 S_1 和与其相邻的两区域 S_0 与 S_2 中, f_4 、 f_5 和 f_0 分别在各子区域中为查询 q 的2近邻。它们满足不等式 $dist(q, f_0) > dist(q, f_4) > dist(q, f_5)$ 。设 $R_1 = dist(q, f_0)$,在区域 S_1 中以 R_1 为半径做扇形 SEC_1 , SEC_1 定义了过滤空间,如图中阴影部分。被阴影部分完全覆盖的网格节点可以直接被过滤;被部分覆盖的网格节点将访问所包含的网格节点或设施,并将位于过滤空间中的网格节点或设施进行过滤。这样一来,算法无需要遍历所有设施便可批量过滤掉不影响候选用户的设施。

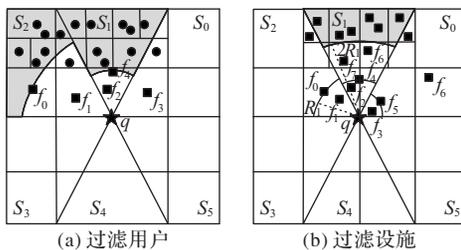


图3 空间属性过滤用户和设施

Fig. 3 Spatial attributes filtering users and facilities

基于Six-region-optimize的过滤算法见算法1。利用两个队列分别缓存用户节点和设施节点如算法1第2)行。对用户进行过滤,不能被过滤的用户存储到候选用户集合 R 中如算法1的第2)~7)行。对设施进行过滤,不能被过滤的设施存储到验证集合 V 中如算法1中第8)~13)行。

算法1 基于Six-region-optimize的过滤算法。

输入 查询设施 q , k 值,用户树 $u-root$ 和设施树 $f-root$,每个区域扇形半径 r ;

输出 集合 R ,集合 V 。

- 1) 初始化 $R \leftarrow \emptyset$, $V \leftarrow \emptyset$;
- 2) Enqueue($u-root$) to U , Enqueue($f-root$) to F ;
- 3) while $U \neq \emptyset$ do
- 4) $u \leftarrow dequeue(U)$;
- 5) $i \leftarrow region(u)$;

- 6) if $dist(u, q) < r[i]$ and $u.k \cap q.k \neq \emptyset$;
- 7) $R \leftarrow u$;
- 8) while $F \neq \emptyset$ do
- 9) $f \leftarrow dequeue(F)$;
- 10) $i \leftarrow region(f)$;
- 11) $max[i] \leftarrow i$ 及相邻区域扇形最大半径;
- 12) if $dist(f, q) < 2 * max[i]$
- 13) $V \leftarrow f$;
- 14) return R, V ;

接下来,介绍Six-region-optimize算法过滤阶段的时间复杂度。该阶段首先需要访问设施KMG-Tree索引进行空间划分,划分需要的时间代价为 $O(\log|V|)$ 。然后,再同时访问用户和设施KMG-Tree索引,根据划分空间对用户和设施进行过滤,该阶段的时间代价为 $O(\log|V| + \log|U|)$ 。整个过滤阶段的总代价为 $O(2\log|V| + \log|U|)$ 。

3.4 基于Six-region-optimize的验证算法

Six-region-optimize的验证算法通过过滤得到候选集 R 中的每个用户 u 进行验证,如果用户 u 同时满足条件1)~3)时将用户存储到最终的结果集 $Result$ 中。

- 1) 用户 u 与查询设施 q 的关键字相似度 $Sim(q, u) \geq \varphi_c$ 。
- 2) 以用户 u 为圆心、 $dist(q, u)$ 为半径做一个圆,遍历集合 V 中的设施,圆中的设施数量小于 k 。
- 3) $dist(u, q)$ 与 $dist(u, NN(u))$ 的比值 $Ratio(q, u) \leq \varphi_r$ 。

以 $\varphi_c = 0.3$ 、 $\varphi_r = 2$ 和 $k = 2$ 为例说明一下。首先,计算候选用户 $u_i(i=0,1,\dots,3)$ 与查询设施 q 的关键字相似度 $Sim(q, u_i)$,如表2所示, $Sim(q, u_0) = 0.71$ 、 $Sim(q, u_1) = 0.63$ 、 $Sim(q, u_2) = 0.09$ 和 $Sim(q, u_3) = 0.67$,所以 u_2 不是查询结果。然后,检查满足关键字相似度要求的候选用户的 k NN是否包含查询设施 q ,例如如图4所示,以 u_1 为中心、 $dist(q, u_1)$ 为半径的圆中有两个设施 f_1 和 f_6 ,所以 u_1 不是查询结果。 u_0 和 u_3 为 q 的R2NN。最后,计算 $dist(q, u_0)$ 与 $dist(u_0, NN(u_0))$ 比值 $Ratio(q, u_0) = 0.7$ 和 $dist(q, u_3)$ 与 $dist(u_3, NN(u_3))$ 比值 $Ratio(q, u_3) = 2.5$ 。 u_3 不满足条件, u_0 为查询结果。

表2 关键字相似度和最近邻距离比值

Tab. 2 Keyword similarity and nearest neighbor distance ratio

用户	$Sim(q, u_i)$	$Ratio(q, u_i)$
u_0	0.71	1.3
u_1	0.63	3.2
u_2	0.09	1.0
u_3	0.67	2.5

基于Six-region-optimize的验证算法见算法2。访问候选用户集合 R 如算法2中的第2)~3)行。验证关键字相似度是否满足查询条件如算法2的第4)行。验证用户的 $RkNN$ 是否包含查询设施 q ,并计算是否满足距离约束如算法2的第5)~11)行。

算法2 基于Six-region-optimize的验证算法。

输入 集合 R 和集合 V ,查询设施 q , k 值,阈值 φ_c 和 φ_r ;

输出 集合 $Result$ 。

- 1) 初始化 $Result \leftarrow \emptyset$;
- 2) foreach $u_i \in R$ do
- 3) $u \leftarrow u_i$;

- 4) if $Sim(u, q) \geq \varphi_k$;
- 5) $temp \leftarrow \emptyset$;
- 6) foreach $f_i \in V$ do
- 7) if $dist(u, f) \leq dist(u, q)$
- 8) $temp++$;
- 9) if $temp < k$ and $Ratio(u, q) < \varphi_r$,
- 10) $Result \leftarrow u$;
- 11) return $Result$;

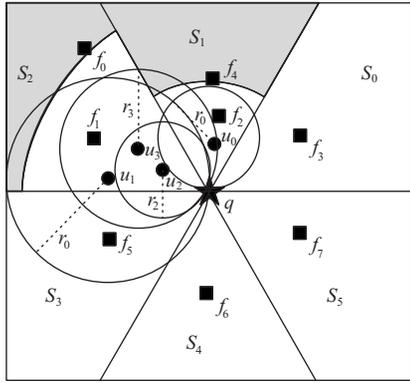


图 4 空间范围验证用户

Fig. 4 Spatial scope authenticating users

接下来,介绍 Six-region-optimize 算法验证阶段的时间复杂度。在算法的过滤阶段已经得到了候选集 R 和验证集 V 。假设候选集的大小为 K ,验证集的大小为 $2K$ 。可以明确的是 $K \ll U$,并且 $2K \ll V$ 。验证阶段需要或每个候选用户发出一个布尔范围查询,每个布尔范围查询大概需要访问 $2/3$ 的验证集中的设施,所以该阶段的时间复杂度为 $O\left(\frac{2}{3}K^2\right)$ 。

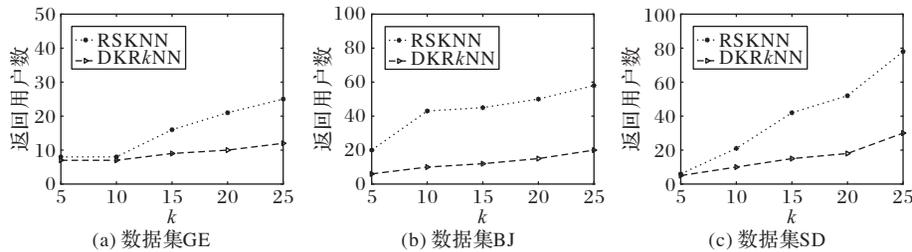


图 5 不同 k 值时 Ranked-RSKNN 和 DKRkNN 返回用户数量

Fig. 5 Number of returned users by Ranked-RSKNN and DKRkNN with different k values

通过图 5 可以发现,在不同数据集上, k 值越大两个算法返回对象数量差距越大,表明 DKRkNN 查询过滤掉的无效的对象越多,从而验证了 DKRkNN 查询返回结果的有效性。

4.2.2 Six-region-optimize 算法的效率分析

通过评估 Six-region-optimize 算法、baseline 算法和 Six-region 算法在数据集大小不同时 CPU 处理时间的消耗分析算法的效率。数据集的大小不同时实验结果如图 6 所示。由图 6 可以看出,在德国(GE)数据集的大小为 6 105 时每个算法在时间上的消耗都非常低, Six-region 和 Six-region-optimize 算法表现得都不是很好。随着数据集的大小不断增加,例如在北京(BJ)数据集的大小为 312 349 和模拟(SD)数据集的大小为 600 000 时,很明显 baseline 算法在时间上的消耗明显增加,而 Six-region 和 Six-region-optimize 算法时耗增加都不是很明显,而 Six-region-optimize 算法时间消耗是最小的,验证了

4 实验与结果分析

4.1 实验准备

实验的配置环境为:64 位 Windows 10、操作系统,Java 编程语言,实验内存为 8 GB,CPU 为 Intel Core i5-4460 处理器。采用的数据集如表 3 所示,包括两个真实数据集德国(GE)的部分地理位置信息和北京(BJ)的部分地理信息以及一个模拟地理位置的数据集(SD)。

4.2 结果分析

首先,对本文提出的 DKRkNN 查询与 Ranked-RSKNN 查询的返回结果进行比较验证该查询返回结果的有效性。然后,使用不同大小的数据集对本文提出的基于 Six-region-optimize 的过滤验证算法进行实验分析,将它与最基本的简称 baseline 的算法和 Stanoi 等^[6]提出的 Six-region 算法进行对比,验证算法高效性。最后,通过改变参数的大小验证本文所提算法的稳定性。

表 3 实验数据集统计信息

Tab. 3 Statistical information of experimental datasets

数据集	实例数	数据集大小
德国(GE)	6 105	201 000
北京(BJ)	312 349	29 000 000
模拟数据(SD)	600 000	59 000 000

4.2.1 DKRkNN 查询的有效性分析

通过比较 DKRkNN 查询和 Ranked-RSKNN 查询在三个数据集上的返回结果数量分析 DKRkNN 查询结果的有效性。如图 5 所示,设阈值 $\varphi_r = 2$ 和 $\varphi_k = 0.5$,设置不同的 k 值为 $\{5, 10, 15, 20, 25\}$ 。

Six-region-optimize 算法的高效性。

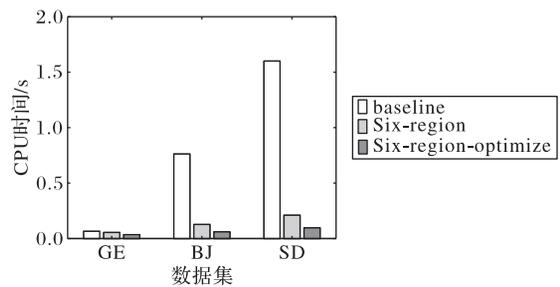


图 6 数据集大小对不同算法时间消耗的影响

Fig. 6 Influence of dataset size on time consumption of different algorithms

4.2.3 Six-region-optimize 算法的稳定性分析

接下来,通过改变参数的值比较 Six-region-optimize 算法、baseline 算法和 Six-region 算法在三个不同数据集上的变化,

分析各个参数对于每个算法的影响。

1) k 值对各算法的影响。

该组实验设阈值 $\varphi_r = 2$ 和阈值 $\varphi_k = 0.5$, 设置不同的 k 值为 $\{5, 10, 15, 20, 25\}$, 三个算法分别在数据集 GE、BJ 和 SD 上的变化情况如图 7 所示。由图 7 可以看出, 随着 k 值增加三个算法在 CPU 时间上的消耗都在增加, 这是因为 k 值越大满足条件的用户越多, 从而过滤掉的数据就越少, 需要验证的数据就越多。但是在数据集较大时, k 值对 Six-region 和 Six-region-optimize 算法的影响很小, 表明了 k 值对这两个算法影响很小, 验证了算法的稳定性。

2) 阈值 φ_r 对各算法的影响。

该组实验设 $k = 5$ 和阈值 $\varphi_k = 0.5$, 设置不同的 φ_r 值为 $\{1.5, 2.0, 2.5, 3.0, 3.5\}$, 三个算法分别在数据集 GE、BJ 和 SD 上的变化情况如图 8 所示。由图 8(a) 可以看出, 三个算法在数据集 GE 上都基本不受影响, 这是因为数据集 GE 非常小, 距离阈值 φ_r 的变化对验证过程影响非常小, 从而三个算法在

CPU 时间上的消耗也比较稳定。由图 8(b)~(c) 可以看出, baseline 算法和 Six-region 算法随着阈值 φ_r 的增加在 CPU 时间上的消耗也有较明显增加, 这是因为这两个算法在过滤阶段没有将大量设施过滤掉, 这样使阈值 φ_r 增加时需要访问的设施也会大量增加。Six-region-optimize 算法不会存在这样的情况, 该算法在过滤阶段已经把大量不影响结果的设施给过滤掉了, 表明了 Six-region-optimize 算法在 CPU 时间消耗上的稳定性。

3) 阈值 φ_k 对各算法的影响。

该组实验设 $k = 5$ 和阈值 $\varphi_d = 9000$, 设置 φ_k 值为 $(0.2, 0.3, 0.4, 0.5, 0.6)$, 三个算法分别在数据集 GE、BJ 和 SD 上的 CPU 时间消耗变化情况如图 9 所示。

由图 9 可以看出, 随着 φ_k 值的增加, 三个算法的 CPU 时间消耗基本无变化, 这是因为关键字相似度只需要在验证阶段计算一次。所以 φ_k 对三个算法都没什么影响, 表明三个算法都是非常稳定的。

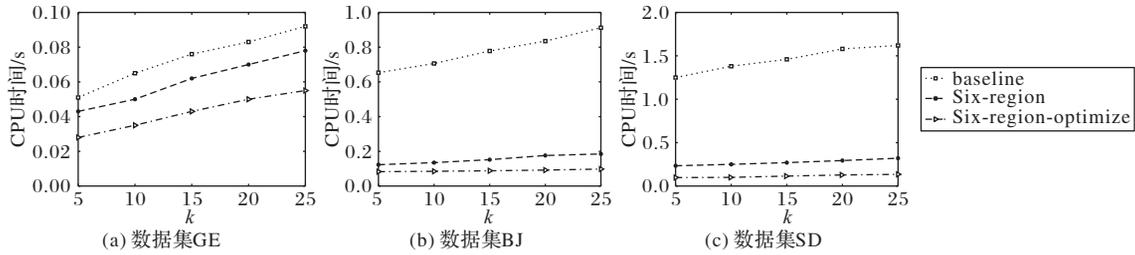


图 7 k 值对不同算法时间消耗的影响

Fig. 7 Influence of k value on time consumption of different algorithms

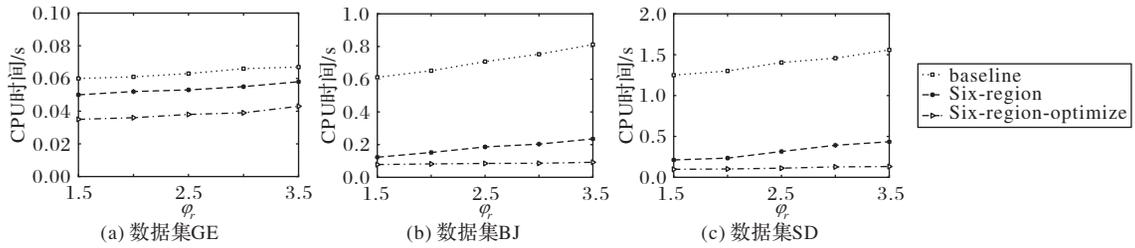


图 8 阈值 φ_r 对不同算法时间消耗的影响

Fig. 8 Influence of threshold φ_r on time consumption of different algorithms

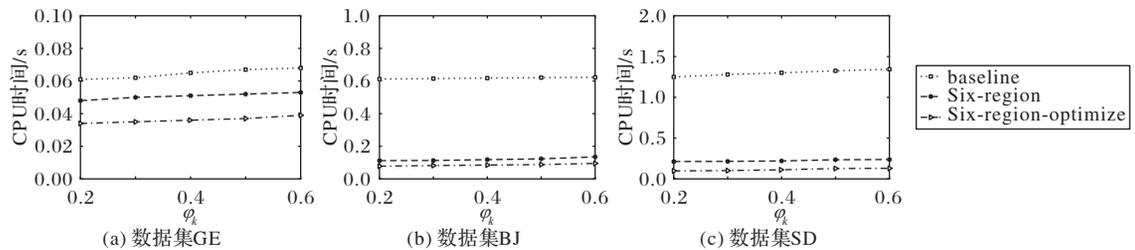


图 9 阈值 φ_k 对不同算法时间消耗的影响

Fig. 9 Influence of threshold φ_k on time consumption of algorithms

5 结语

本文研究了距离-关键字相似度约束的双色反 k 近邻查询(DKR k NN)问题, 并提出了一种基于多分辨率网格树索引查询处理框架。该框架根据区域划分规定过滤区域, 访问多分辨率网格树, 得到用户候选集和影响用户成为查询结果的设施, 对候选用户进行验证, 得到最终的结果集。基于两个真实数据集和一个模拟数据集进行实验, 实验结果验证了本文

算法的有效性和高效性, 表明所提框架能够有效处理该查询。但是随着社会的发展, 查询的需求会越来越复杂, 空间文本信息可能不再是以关键字的方式存在, 有可能是一段文字。因此, 未来的研究需要对空间文本做语义分析从而为用户提供更高质量的查询体验。

参考文献 (References)

[1] LU J, LU Y, CONG G. Reverse spatial and textual k nearest

- neighbor search [C]// Proceedings of the 2011 ACM IGMOD International Conference on Management of Data. New York: ACM, 2011: 349-360.
- [2] ZHAO P, FANG H, SHENG V S, et al. Monochromatic and bichromatic ranked reverse boolean spatial keyword nearest neighbors search [J]. World Wide Web, 2017, 20(1): 39-59.
- [3] XIA S, ZHANG Z, LI W, et al. GBNRS: a novel rough set algorithm for fast adaptive attribute reduction in classification [J]. IEEE Transactions on Knowledge and Data Engineering, 2020 (Early Access): 1-1.
- [4] XIA S, PENG D, MENG D, et al. A fast adaptive k -means with no bounds [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence. 2020(Early Access): 1-1.
- [5] XIA S, WANG G, CHEN Z, et al. Complete random forest based class noise filtering learning for improving the generalizability of classifiers [J]. IEEE Transactions on Knowledge and Data Engineering, 2019, 31(11): 2063-2078.
- [6] STANOI I, AGRAWAL D, ABBADI A E. Reverse nearest neighbor queries for dynamic databases [C]// Proceedings of the 2000 ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery. New York: ACM, 2000: 44-53.
- [7] TAO Y, PAPADIAS D, LIAN X. Reverse k NN search in arbitrary dimensionality [C]// Proceedings of the 2004 30th International Conference on Very Large Data Base. Amsterdam: Elsevier, 2004: 744-755.
- [8] TAO Y, PAPADIAS D, LIAN X, et al. Multidimensional reverse k NN search [J]. The VLDB Journal, 2007, 16(3): 293-316.
- [9] CHEEMA M A, LIN X, ZHANG Y, et al. Lazy updates: an efficient technique to continuously monitoring reverse k NN [J]. Proceedings of the VLDB Endowment, 2009, 2(1): 1138-1149.
- [10] CHEEMA M A, SHEN Z, LIN X, et al. A unified framework for efficiently processing ranking related queries [C]// Proceedings of the 2014 17th International Conference on Extending Database Technology. Berlin: Springer, 2014:427-438.
- [11] CHEEMA M A, ZHANG W, LIN X, et al. Efficiently processing snapshot and continuous reverse k nearest neighbors queries [J]. The VLDB Journal, 2012, 21(5): 703-728.
- [12] LEE K C K, ZHENG B, LEE W C. Ranked reverse nearest neighbor search [J]. IEEE Transactions on Knowledge and Data Engineering. 2008, 20(7): 894-910.
- [13] YIU M L, PAPADIAS D, MAMOUI N, et al. Reverse nearest neighbors in large graphs [J]. IEEE Transactions on Knowledge and Data Engineering, 2006, 18(4): 540-553.
- [14] 李佳佳,李雨现,夏秀峰,等. 面向时间依赖路网的连续 k 近邻查询[J]. 计算机科学与探索, 2019, 13(5): 788-799. (LI J J, LI Y X, XIA X F, et al. Continuous k -neighbor query for time dependent road network [J]. Journal of Frontiers of Computer Science and Technology, 2019, 13(5): 788-799.)
- [15] HIDAYAT A, YANG S, CHEEMA M A, et al. Reverse approximate nearest neighbor queries [J]. IEEE Transactions on Knowledge and Data Engineering, 2018, 30(2): 339-352.

This work is partially supported by the National Natural Science Foundation of China (61702344).

ZHANG Hao, born in 1996, M. S. candidate. His research interests include big data.

ZHU Rui, born in 1982, Ph. D., associate professor. His research interests include big data.

SONG Fuyao, born in 1995, M. S. candidate. His research interests include big data.

FANG Peng, born in 1994, M. S. candidate. His research interests include big data of new media.

XIA Xiufeng, born in 1964, Ph. D., professor. His research interests include database, data warehouse, data mining.