

Top-K Ranking Spatial Queries over Filtering Data

¹Lakkapragada Prasanth, ²Krishna Chaitanya,
¹Student, ²Assistant Professor,
 NRL Institute of Technology, Vijayawada

Abstract: A spatial preference query ranks objects based on the qualities of features in their spatial neighborhood. For example, using a real estate agency database of flats for lease, a customer may want to rank the flats with respect to the appropriateness of their location, defined after aggregating the qualities of other features (e.g., restaurants, cafes, hospital, market, etc.) within their spatial neighborhood. Such a neighborhood concept can be specified by the user via different functions. But they are accessing services only non filtering methodology. In this paper we propose to develop Top-k preferable product categorization can be defined efficient process generation. Our experimental results show efficient filtering rules with semantic data representation.

Keywords: *Query processing, spatial databases, top-k spatial preference query.*

I. INTRODUCTION

Spatial database systems manage large collections of geographic entities, which apart from spatial attributes contain non-spatial information (e.g., name, size, type, price, etc.). In this paper, we study an interesting type of preference queries, which select the best spatial location with respect to the quality of facilities in its spatial neighborhood. Given a set D of interesting objects (e.g., candidate locations), a top-k spatial preference query retrieves the k objects in D with the highest scores. The score of an object is defined by the quality of features (e.g., facilities or

services) in its spatial neighborhood. As a motivating example, consider a real estate agency office that holds a database with available flats for lease. Here “feature” refers to a class of objects in a spatial map such as specific facilities or services.

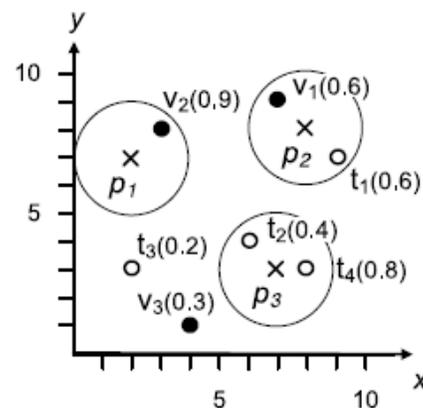


Figure 1: Spatial area containing data and feature objects.

With the popularization of geotagging information, there has been an increasing number of Web information systems specialized in providing interesting results through location-based queries. However, most of the existing systems are limited to plain spatial queries that return the objects present in a given region of the space.

For example, Figure 1 presents a spatial area containing data objects p (hotels) together with feature objects t (restaurants) and v (cafes) with their respective scores (e.g. rating). Consider a tourist interested in hotels with good restaurants and cafes

in their spatial neighborhood. The tourist specifies a spatial constraint (in the figure depicted as a range around each hotel) to restrict the distance of the eligible feature objects for each hotel. Thus, if the tourist

wants to rank the hotels based on the score of restaurants, the top-1 hotel is p3(0.8) whose score 0.8 is determined by t4. However, if the tourist wants to rank the hotels based on cafés, the top-1 hotel is p1(0.9) determined by v2.

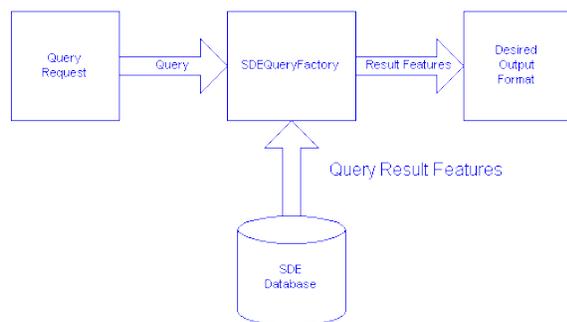


Figure 2: Interoperability of the query processing in spatial query processing.

In this paper, we propose a novel approach for processing spatial preference queries efficiently. The main difference compared to traditional top-k queries is that the score of each data object is defined by the feature objects that satisfy a spatial constraint (for example range constraint). Therefore, pairs of data and feature objects need to be examined to determine the score of an object. Our approach relies on mapping of pairs of data and feature objects to a distance-score space, which in turn allows us to identify the minimal subset of pairs that is sufficient to answer all spatial preference queries. Capitalizing on the materialization of this subset of pairs, we present an efficient algorithm that improves query

processing performance by avoiding examining the spatial neighborhood of data objects during query execution.

II. RELATED WORK

Several approaches have been proposed for ranking spatial data objects. The reverse nearest neighbor (RNN) query was first proposed by Korn and Muthukrishnan. Then, Xia *et al.* studied the problem of retrieving the top-k most influential spatial objects, where the score of each spatial data object p is defined as the sum of the scores of all feature objects that have p as their nearest neighbor. Yang *et al.* studied the problem of finding an optimal location.

The main difference compared to is that the optimal location can be any point in the data space and not necessarily an object of the dataset, while the score is computed in a similar way to. The aforementioned approaches define the score of a spatial data object p based on the scores of feature objects that have p as their nearest neighbor and are limited to a single feature set. Differently, Yiu *et al.* first considered computing the score of a data object p based on feature objects in its spatial neighborhood from multiple feature sets. To this end, three different spatial scores were defined: range, nearest neighbor, and influence score; and different algorithms were developed to compute top-k spatial preference queries for these scores.

To solve this type of queries, many of the researchers proposed several methods. In 2004, F. Ilyas *et al.* introduced a new rank-join algorithm that made use of the individual orders of its inputs to produce join results ordered on a user-specified scoring function. They had experimentally evaluated their proposed rank join operators and analyze its performance. In 2006, Rank-aware query

optimization framework by Ihab F. Ilyas et.al. fully integrated the rank-join operators into relational query engines and shown the performance of the proposed framework. In 2007 Yiu et.al., proposed the Branch and bound (BB) and Feature join algorithm (FJ) that rank objects based on the qualities of features. They proved that their proposed work is better than simple and Group probing algorithms with real and synthetic data. The top-k queries produce ordered result by using some calculated score. Generally, users are interested in top-k join result. For this, the top-k queries require joins to produce top-k result. The relational processors should not process the ranking queries with join efficiently.

III. BACKGROUND WORK

3.1. Definitions and Index Structures

Let F_c be a feature data set, in which each feature object $s \in F_c$ is associated with a quality $!sP$ and a spatial point.

We assume that the domain of $!sP$ is the interval $1/20; 1_.$ As an example, the quality $!sP$ of a restaurant s can be

obtained from a ratings provider. Let D be an object data set, where each object $p \in D$ is a spatial point. In other words, D is the set of interesting points (e.g., hotel locations) considered by the user.

Probing Algorithms

We first introduce a brute-force solution that computes the score of every point $p \in D$ in order to obtain the query results. Then, we propose a group evaluation technique that computes the scores of multiple points concurrently.

```

Algorithm 1. Simple Probing Algorithm
algorithm SP(Node N)
1: for each entry  $e \in N$  do
2: If N is nonleaf then
3: read the child node  $N_0$  pointed by e;
4: SP( $N_0$ );
5: else
6: for  $c : 1$  to  $m$  do
7: If  $!p\delta eP > _$  then . upper bound score
8: compute  $_c\delta eP$  using tree  $F_c$ ; update
 $_p\delta eP$ ;
9: If  $!eP > _$  then
10: update  $W_k$  (and  $_$ ) by e;

```

Algorithm 1 is a pseudo code of the simple probing (SP) algorithm, which retrieves the query results by computing the score of every object point. The algorithm uses two global variables: W_k is a min-heap for managing the top-k results and represents the top-k score so far (i.e., lowest score in W_k). Initially, the algorithm is invoked at the root node of the object tree (i.e., $N \neq D$:root). The procedure is recursively applied (at Line 4) on tree nodes until a leaf node is accessed.

Upper Bound Score Computation

It remains to clarify how the (upper bound) scores $T_c\delta eP$ of nonleaf entries (within the same node N) can be computed concurrently. Our goal is to compute these upper bound scores such that . the bonds are computed with low I/O cost, and . the bonds are reasonably tight, in order to facilitate effective pruning. To achieve this, we utilize only level-1 entries (i.e., lowest level nonleaf entries) in F_c for deriving upper bound scores because: 1) there are much fewer level-1 entries than leaf entries (i.e.,

points), and 2) high-level entries in F_c cannot provide tight bounds.

IV. PROPOSED WORK

A Spatial preference query, ranks the spatial objects based on quality of its neighbor facilities. For instance a tourist might retrieve a sorted list of hotels based on the facilities around that (e.g. restaurant, hospital, market, etc.). Assume that p is our point of interest (e.g. a hotel) and we have m type of facilities (e.g. restaurant means $m=1$ and park means $m=2$). Then assume that $n m f$ is n -th facility from type m (e.g. Restaurant A). First we retrieve a list of candidates for P according to Table 1. Table 1 shows how one of the methods choose the primary candidates.

Table 1 Candidate Selection Criteria

Method	
Nearest Neighbor	$\min(d(p, f_m^n))$
Range Score	$d(p, f_m^n) < R$
Influence Score	All

As we can see, Nearest Neighbor, from each type m retrieves n -th element of that $(n m f)$ which has the minimum distance with p . Range score retrieves a list of items which have at least distance (d) of pre-defined R with P . Influence score retrieves all the items for further computation. Afterwards, We define Score of point P according to the following equation:

$$C_i S = \sum Agg w \times \alpha (1)$$

Where, Agg denotes the aggregation function which can be maximum or sum. w is equal to the weight or quality of item (e.g. hotel with 5 star can have weight

of 5 and hotel with one star can have weight of 1) and i is an index of retrieved candidates. α is influence function which is equal to 1 for Nearest Neighbor and Range score and is equal to the equation 2 for Influence score. $() R f p d i m , 2 - \alpha = (2)$ Where d denotes the distance between point P and facility i of category m . and R is a pre-defined radius. Then the result of Top-K spatial preference query is a sorted list of S_p for all point of interests (P).

V. EXPERIMENTAL RESULTS

In this section, we evaluate our proposed algorithm (SFA) and we compare SFA against the algorithms developed by Yiu *et al.* denoted as GP , BB , BB^* , and FJ . All algorithms were implemented in Java and executed on a PC with 3GHz Dual Core AMD Processor with 2GB RAM. The datasets were indexed by an R-tree (aR-tree for [16, 17]) with block size of 4KB. We used an LRU memory buffer with a fixed size of 0.2% of the size of the total number of objects stored in O and F_i . We report the average values of 20 experiments, and in each experiment we recreate all datasets and indexes to factor out the effects of randomization. In all experiments, we measured the total execution time (referred to as response time) and number of I/Os. All charts are plotted using a logarithmic scale on the y-axis.

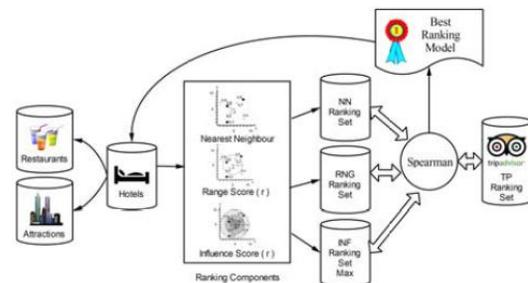


Figure 3: Experimental setup progression.

Experimental Settings

We conduct experiments using both synthetic and real datasets. First, we perform experiments using uniform distribution (UN) for the spatial locations of data and feature objects and for the score of the feature objects (within the range [0, 1]). We also generate a synthetic dataset (CN) that resembles the real world: (1) there exist multiple city centers (centroids) with higher occurrences of data objects, (2) there exists a higher probability of finding feature objects nearby the city centers (centroids). Appendix C.1 provides a detailed description of CN including a plot of a generated dataset. We use the synthetic dataset (CN) as our default dataset. By default, the non-spatial score of the feature objects is a uniformly generated value within the range [0, 1].

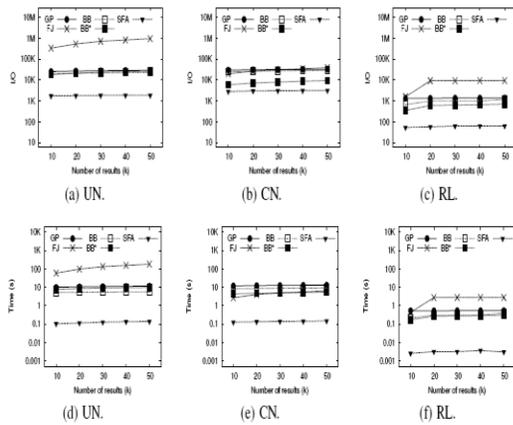


Figure 4: **Effect of different data distributions {UN,CN,RL} on I/O and response time (range score).**

Query Processing Performance

Range Score. In Figure 4, we use our default setup and study the number of I/Os and the response time for all datasets, while varying k . Figure 4(a) presents the I/O cost using the UN dataset. The performance of GP is stable because it always computes the score

of all data objects. FJ requires a much higher number of I/Os, as it needs to access many leaf entries of the feature R-trees in order to report the correct top- k result set. The branch-and-bound algorithms (BB and BB*) perform slightly better than GP for this setup. However, SFA results in one order of magnitude fewer I/Os than the best of its competitors. In Figure 4(b), we plot the number of I/Os for the CN dataset. BB* performs better than GP, BB, and FJ due to the employed pruning. However, SFA reduces even further the number of required I/Os compared to BB* and scales better than BB* for increasing value of k . In Figure 4(c), the I/O cost for the real dataset (RL) is presented. Again, SFA outperforms all other algorithms (in terms of I/Os) by at least one order of magnitude. This experiment indicates that SFA performs efficiently for a wide range of different datasets. Figures 4(d), 4(e) and 4(f) depict the response time for the same experimental setups respectively

VI. CONCLUSION

In this paper, we present a novel approach for boosting the performance of top- k spatial preference query processing. At the heart of our framework lies a mapping of pairs of data and feature objects to a distance-score space, which enables us to identify the minimal subset of pairs necessary to answer any ranked spatial preference query. By materializing this subset of pairs, we present efficient algorithms for query processing that result in improved performance. Furthermore, we describe an efficient algorithm for materialization and elaborate on useful properties that reduce the cost of maintenance. Our experimental evaluation demonstrates that our approach reduces I/Os and response time by more

than one order of magnitude compared to the state-of-the-art algorithms in most of the setups.

VII. REFERENCES

- [1] C. Böhm, B. C. Ooi, C. Plant, and Y. Yan. Efficiently processing continuous k-nn queries on data streams. In *Proc. of Int. Conf. on Data Engineering (ICDE)*, pages 156–165, 2007.
- [2] S. Börzsönyi, D. Kossmann, and K. Stocker. The skyline operator. In *Proc. of Int. Conf. on Data Engineering (ICDE)*, page 421430, 2001.
- [3] S. Chaudhuri, N. N. Dalvi, and R. Kaushik. Robust cardinality and cost estimation for skyline operator. In *Proc. of Int. Conf. on Data Engineering (ICDE)*, page 64, 2006.
- [4] Y. Du, D. Zhang, and T. Xia. The Optimal-Location query. In *Proc. of the Int. Symposium on Spatial and Temporal Databases (SSTD)*, pages 163–180, 2005.
- [5] <https://developers.google.com/maps/>
- [6] <http://developer.yahoo.com/maps/>
- [7] C. Spearman, The Proof and Measurement of Association between Two Things, *The American Journal of Psychology*, Vol. 15, No. 1 (Jan., 1904), pp. 72-101.
- [8]. Yiu. M. L, Dai.X,Mamoulis.N, Vaitis.M. “Ranking Spatial data by quality preferences”.IEEE Transaction on Knowledge and Data Engineering,Vol.23, 3, pp.433 – 446, 2011.
- [9]. Yong Zhang, Lizhu Zhou, Jun Chen. “Layered r-tree: an index structure for three dimensional points and lines”. International Archives of Photogrammetry and Remote Sensing, Vol. XXXIII,B4.
- [10]. Antonio Leopoldo, Corral Liria. Algorithms for the Processing of Spatial Queries using R-trees. The

Closest Pairs Query and its Application on Spatial Databases Almeria, January 2002

- [11]. Yunjunga gao,“Optimal-location-selection query processing in spatial database”knowledge and data engineering, august 2009.